

OXFORD
UNIVERSITY PRESS

Briefings in Bioinformatics

Computational analysis of protein interaction networks for infectious diseases

Journal:	<i>Briefings in Bioinformatics</i>
Manuscript ID:	BIB-15-0058.R2
Manuscript Type:	Paper
Date Submitted by the Author:	08-Jul-2015
Complete List of Authors:	Pan, Archana; Pondicherry University, Centre for Bioinformatics Lahiri, Chandrajit; Sunway University, Department of Biological Sciences Rajendiran, Anjana; Pondicherry University, Centre for Bioinformatics Shanmugham, Buvaneswari; Pondicherry University, Centre for Bioinformatics
Keywords:	Protein interaction network, Infectious disease, Pathogen, Computational analyses, Centrality, Modularity

SCHOLARONE™
Manuscripts

Preview

1
2
3
4 ***Computational analysis of protein interaction networks for infectious diseases***
5
6
7
8
9
10
11
12
13

14 **Authors:** Archana Pan^{1#*}, Chandrajit Lahiri^{2,#}, Anjana Rajendiran¹ and
15
16 Buvaneswari Shanmugham¹
17
18
19
20

21 ¹Center for *Bioinformatics*, School of Life Sciences, Pondicherry University,
22
23 Puducherry – 605014, India
24

25 ²Department of Biological Sciences, Sunway University,
26
27 47500 Bandar Sunway, Selangor, Malaysia
28
29
30
31
32

33 #These authors contributed equally to this work
34
35
36
37

38 ***Corresponding Author:** Archana Pan, **E-mail :** archana@bicpu.edu.in
39

40 **Phone:** +91 413 2654584
41

42 **Fax:** +91 413 2655211
43
44
45
46
47

48 **Authors:** Archana Pan, **E-mail :** archana@bicpu.edu.in
49

50 Chandrajit Lahiri, E-mail : chandrajithlahiri@gmail.com
51

52 Anjana Rajendiran, Email: anjana@mails.bicpu.edu.in
53

54 Buvaneswari Shanmugham, Email: buvanisuriya@bicpu.edu.in
55
56
57
58
59
60

Abstract

Infectious diseases caused by pathogens, including viruses, bacteria and parasites, pose a serious threat to human health worldwide. Frequent changes in the pattern of infection mechanisms and the emergence of multidrug resistant strains among pathogens have weakened the current treatment regimen. This necessitates the development of new therapeutic interventions to prevent and control such diseases. To cater to the need, analysis of protein interaction networks (PINs) has gained importance as one of the promising strategies. The present review aims to discuss various computational approaches to analyse the PINs in context to infectious diseases. Topology and modularity analysis of the network with their biological relevance, and the scenario till date about host-pathogen and intra-pathogenic protein interaction studies were delineated. This would provide useful insights to the research community thereby enabling them to design novel biomedicine against such infectious diseases.

Keywords: Protein interaction network, Infectious disease, Pathogen, Computational analyses, Centrality, Modularity.

Author Biography

Archana Pan, PhD in Bioinformatics, currently serves as an Assistant Professor at Centre for Bioinformatics, Pondicherry University, Pondicherry, India. Her research interests include comparative genomics, molecular evolution and drug design.

Chandrajit Lahiri, PhD in Molecular Microbiology, is currently a Senior Lecturer at Department of Biological Sciences, Sunway University, Selangor, Malaysia. His research interests encompass Systems Biology and Evolutionary Bioinformatics.

Anjana Rajendiran is registered for PhD at Centre for Bioinformatics, Pondicherry University, Pondicherry, India. Her research interests include comparative genomics, machine learning, and microRNA prediction.

Buveneswari Shanmugham is registered for PhD at Centre for Bioinformatics, Pondicherry University, Pondicherry, India. Her research interests lie in comparative genomics, structural bioinformatics and genetic algorithms.

Key Points

1. Infectious diseases have posed serious health concerns worldwide
2. Conventional approaches have become almost ineffective in dealing with the issue
3. Non-conventional computational approach entailing protein interaction network analysis has gained importance to give meaningful directions
4. Topological and Modularity analyses of PINs can be employed by researchers to obtain essential proteins as key therapeutic targets
5. Analyses involving these would pave the way to succeed in generating novel biomedicines

Introduction

Infectious diseases have been threatening human population since time immemorial. These have become ever-increasing worldwide public health concern, with parasitic, bacterial and viral diseases, representing more than half of the leading causes of morbidity and mortality. While viral influenza vaccines are to be reformulated annually, several other viral infectious diseases, such as those from hepatitis C and HIV-1 are a cause of panic since decades. Parasitic diseases like malaria and multidrug resistant bacterial strains of *Mycobacterium* and *Salmonella* are on-going pandemics. Different emerging infectious diseases viz., nosocomial infections caused by *Acinetobacter*, swine H1N1 influenza, avian H5N1 influenza, severe acute respiratory syndrome (SARS), and dengue fever have posed themselves to be new constant threats [1].

To deal with these severe pathogenic threats, several health intervention strategies have been undertaken. However, the prospects for finding new vaccines or antibiotics against such pathogens are especially poor. This is due to the ever-changing mechanism of infecting the host as in the cases of viruses [1]. It might also be due to the blockades provided by the outer membrane to the entry of some existing antibiotics in case of gram negative bacteria [2]. Thus, it is quite evident that the conventional strategies for dealing with such deadly pathogens would be less effective or ineffective completely, to emerge victorious against their strategies to evade therapeutic interventions. In such cases, the complexities posed can be solved by adopting some non-conventional computational approaches.

Over the last few decades, biologists understood gradually that a set of complex interactions between the numerous constituents of a cell, gives rise to different biological phenotypes. Amongst these, proteins, being the functional unit of the cell of any living organism, always act

1
2
3 in unison with others to achieve specific functional goals viz., transcriptional
4 activation/repression; immune, endocrine, and pharmacological signaling; cell-to-cell
5 interactions; and metabolic and developmental control [3]. These protein-protein interactions
6 (PPI) lead to a mosaic mesh or network of interactions, commonly known as protein interaction
7 networks (PINs). Analyses of such PINs are increasingly serving as the non-conventional
8 approach to understand the complexity of infectious diseases. However, the augmentation of the
9 PINs, created from high-throughput experimental and/or computational data, has necessitated
10 effective analytical techniques for those networks, to be used to unravel the molecular basis of
11 the aforementioned infectious diseases. The current review entails different computational
12 approaches for analysing protein interaction networks expected to be involved in the interaction
13 mechanism of infection. These might lead to find avenues for the identification of novel targets
14 and render them as systems biomedicines.

31 32 33 34 **The necessity of the generated PIN**

35
36 With the advances of the post genomic era, there has been an enormous increase in the
37 investigations upon the structure, function and control of the participating proteins as key
38 regulators in diseases. This is because, the identification of a handful of proteins to be targeted is
39 considered as the objective of the whole intervention process. The numbers of proteins, as
40 targets, should always be limited, to improve the efficacy and specificity of a well-defined drug.
41 However, ensuring a limited number of proteins from an array becomes an ever challenging task
42 to the conventional experimentalists. Thus, new approaches, for generating viable candidates as
43 interventional targets for infectious diseases, are need of the hour.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 The complexities of any infectious diseases are mainly due to the intricate interactions between
4 sets of proteins involved in the process. Interactions between proteins are visualized by networks
5 created by mapping those complex interactions. These protein interaction networks (PINs) have
6 gradually gained importance in an attempt to address the complexities of the diseases. Such
7 mapping can be done based on a number of experimental data sources including, but not limited
8 to, two-hybrid systems [4], mass spectrometry [5], protein chip technologies [6]. They can also
9 be generated through various computational approaches encompassing genome-based [7, 8],
10 sequence-based [9, 10], structure-based [11, 12] and machine-learning-based techniques [13, 14].
11 However, analysing these networks, to achieve the ultimate goal of limiting target sets for health
12 intervention, now becomes the most challenging task.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

29 **The resources of PPIs**

30
31 While the high throughput techniques generated interaction data for proteins, initiatives were
32 taken to integrate them and prepare comprehensive open databases for further analyses. There
33 are a number of standardized open sources each having a different style of representing the
34 protein interaction datasets. They are mostly based on the organisms worked upon in detail and
35 of basic interest amongst researchers. Of these, Human Protein Reference Database (HPRD)
36 stores information on human protein interactions, along with protein functions, post-translational
37 modifications (PTMs), enzyme-substrate relationships, and subcellular localization [15]. Sub-
38 categorised HomoMINT [16] arises from the Molecular Interaction Database (MINT) [17] which
39 comprises interactions, inferred from orthologs in model organisms. For the yeast PPI data,
40 special importance has been given in the Biological General Repository for Interaction Database
41 (BioGRID) [18]. The current BioGRID release [May 2015, version 3.3.124] lists 287,619 non-
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

redundant yeast protein interactions thereby making it the largest database for this organism besides more than 30 others, having a total of [574,378](#) non-redundant interactions. Another such database focusing on yeast data is the Database of Interacting Proteins (DIP) integrating data from the correlation of protein sequence and RNA expression profiles through a carefully curated computational process [19].

Besides the above mentioned focused databases, there are others, having listed the protein interactions from a set of organisms. These are the Munich Information Center for Protein Sequences (MIPS) [20], the Biomolecular Interaction Network Database (BIND) [21], a component of the Biomolecular Object Network Databank (BOND), the Search Tool for Recurring Instances of Neighbouring Genes/Proteins (STRING) [22] and IntAct [23], each having its own uniqueness. MIPS lists a description and the binding regions of interacting partners. BIND highlights the interactions between two or more molecules which form functional molecular complex units and pathways arising from those interacting in a sequence. The STRING database entails both physical and functional associations derived from genomic context, high-throughput experiment, coexpression and previous knowledge. Apart from interaction data, IntAct enlists interactions between DNA, RNA, and small-molecules. Furthermore, some databases including STRING [22], GeneMANIA [24], FunCoup [25] and ConsensusPathDB [26] provide a highly comprehensive data by integrating PPIs from other online resources. [STRING imports PPI data from different primary databases, including MINT, HPRD, BioGRID, DIP, BIND, IntAct and PDB.](#) GeneMania [provides](#) functionally similar genes for the query gene [list](#) along with interactive functional association network utilizing information from GEO, BioGRID, Pathway Commons and I2D. FunCoup includes information on functional couplings between genes and gene products based on gold standards (KEGG, Corum, iRefindex

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

etc.). ConsensusPathDB enlists a seamless interaction network from different public resources, including BIND, BIOGRID, DIP, IntAct, MINT and MIPS-MPPI. There are different resources providing intra-pathogenic and host-pathogen PPI data, such as PATRIC, PRIMOS, HPIDB, PHI, VirusMentha, VirusHostNet [27-33]. The number of intra-pathogenic and host-pathogen PPIs is tabulated in Supplementary Table 1.

The technicalities of PIN analysis

The databases stated above list the interactions of the proteins from existing empirical and theoretical results. As such, an attempt to construct a network or an interactome, by integrating those interactions, might yield one, which can be random like the one proposed by Erdős and Renyi [34] or a small-world type proposed by Watts and Strogatz [35]. Both these types build up a fairly homogeneous network in which, each node has approximately the same number of links. However, only those interactomes, which strictly follow the power law, are free of a characteristic scale. In these cases, the connectivity distribution, $P(k)$, of a node in a network getting connected to k other nodes, **decays** exponentially for large values of k . These scale-free networks are essentially the real world networks [36] with a heavy tailed degree distribution. Thus, it is imperative to construct biologically viable real networks, comprising the proteins responsible for the infectious diseases. Their subsequent analyses, in essence, would then lead us to our ultimate goal of identifying important targets for health intervention.

The analyses of the interactomes

An overview of various computational approaches for protein interaction network analysis is illustrated in Figure 1.

Topological analyses to identify an important protein

In order to identify the key proteins in a PIN, the importance of the protein is correlated with the number of its interacting protein partners. This gives rise to the concept of such proteins becoming central to a particular network. This is the most basic concept in terms of biological importance and is defined as the degree centrality (DC) of the protein in a network of interacting proteins. Indeed, high degree proteins (or hubs) are known to correspond to the essential proteins in a network [37]. However, DC is a local and static metric, as it considers only the directly connected neighbours of a protein in a static state. Thus, DC, being the local property of a protein in the network, does not bring out the importance of the protein on a global scale. To indicate such importance based on a protein's global relevance in the network, researchers resort to other centrality measures. These are Closeness centrality (CC), Betweenness centrality (BC) and Eigenvector centrality (EC) [38]. These four important concepts of centrality measures reportedly have been utilized for biological network analyses [39-41].

It is understood that, being the most basic of the centrality measures, DC generally refers to the protein involved in a large number of interactions in a network. However, these interactions might not be in a sequential order so as to carry out particular functions during the primary stages of infection by a pathogen. Conceptually, CC might take care of this fact as it reflects the protein, which is typically "close" to, and can communicate sequentially with the other proteins in the network. Thus, CC is a measure based on the interacting distance of a protein to all other proteins in a network. It is defined as the reciprocal of the total interacting distance from a protein to all other proteins in the network. Again, in a complex phenotype like virulence in infectious diseases, there might be simultaneous interaction of a protein with others to render different functions at the same time. Thus, an important protein should be typically the one,

1
2
3 which lies on a high proportion of interactions between other proteins in the network. An
4 analysis with BC might bring out this fact. Thus, a better measure compared to the DC and CC
5 would be BC, since it would reflect the importance of the protein with respect to its
6 indispensability as it would form the bridge between important hubs of network thereby
7 becoming important. BC of a protein is defined as the number of shortest interacting paths
8 passing through it. However, the ultimate idea of a protein in a network to be important lies in
9 the fact that it should be connected to other important proteins in the network. EC might come
10 into play in such cases [42]. EC brings out the relative importance of the proteins in the network
11 by weighting the connections to other important proteins compared to those of low importance
12 [43].

13
14
15
16
17
18
19
20
21
22
23
24
25
26
27 It has been observed that topological features like DC and BC have gained much importance in
28 serving as attractive drug targets [36, 44, 45]. However, despite their potential to locate such
29 targets, these measures lack in the specificity and/or selectivity along with the high risk of side
30 effects. These, in turn, result in a high likelihood of causing lethality as determined
31 experimentally in the yeast PIN [46, 47]. As lethality is an undesirable attribute in most of the
32 drug discovery applications [3], an alternative measure for betweenness can be thought of. This
33 is known as bridging centrality and proteins with high bridging centrality mainly serve as
34 bottlenecks between two modules. This has been shown to be less lethal, with a value of 34%
35 compared to 42 for BC and 48 for DC in case of yeast PIN [48].

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
There are other topological properties which have been utilized to measure the compactness and
reachability amongst the interacting proteins in the network. One of these is the average path
length (APL) which determines the mean of the lengths of the shortest paths between all protein

1
2
3 components of the network [49-51]. The other is network diameter which measures the longest
4
5 distance between two constituent components [49].
6
7
8
9

10 **Network decomposition to identify set of important proteins**

11
12 In general, the PIN for an infectious disease would be on a large scale. Thus, as discussed in the
13
14 above section, a focus to target just one protein, for therapeutic health intervention, may be of
15
16 less importance. This might necessitate a decomposition of those large networks to a core of
17
18 highly interacting proteins through the k -core analysis approach [52]. This essentially peels off
19
20 the proteins connected at the edges, gradually, until the innermost core is reached. After this
21
22 core, a step further decomposes the network, thereby making this the innermost core with highly
23
24 connected proteins, interacting with each other. Thus, they can be considered to be the most
25
26 important ones [42].
27
28
29
30

31
32 Technically speaking, the k -core of the graph G is obtained by recursively removing all the
33
34 vertices of degree less than k , until all vertices in the remaining graph have at least degree k , by
35
36 which the complex network can be decomposed [53, 54].
37
38
39
40

41 **Modularity analyses and functional annotation of clustered proteins**

42
43 The concept of k -core, as discussed above, is one of the metrics to determine the modularity of a
44
45 network. A modular network groups the components on the basis of their common properties to
46
47 bring out significant underlying principles. Analyses of these networks become increasingly
48
49 useful for PINs. This is due to the biological phenomenon of proteins aggregating into
50
51 complexes, rendering them as functional modules which unify the cohesive components of a
52
53 molecular function. The identification of such highly correlated functional modules of proteins
54
55
56
57
58
59
60

1
2
3 can be done by clustering analyses. These protein modules from one species can then be utilised
4
5 to rationally map and thereby annotate the unannotated proteins in other related genomes.
6
7

8 Besides k -core, the clustering techniques can identify cliques. A clique is an induced complete
9
10 subnetwork where each component vertex is connected to each other. This gives rise to a
11
12 clustering coefficient of 1 for each of the component vertices. Parametric indices like maximum
13
14 clique centrality (MCC), maximum neighbourhood centrality (MNC) and density of maximum
15
16 neighbourhood centrality (DMNC) are offshoots of these concepts and has been utilised lately
17
18 [42]. Such densely connected subnetworks are expected to form functional units to carry out
19
20 unique biological processes.
21
22
23

24 While such density based traditional clustering method is in good practice amongst researchers,
25
26 new approaches through non-traditional methods have started gaining importance. This is
27
28 because of their ability to analyse the modularity of the PPI networks with more accuracy. These
29
30 include the graph-theoretic, topology-based, flow-based, statistical, and domain knowledge-
31
32 based approaches (data fusion, GO integration) besides the distance-based methods [3]. Of these,
33
34 the topology- and distance-based modularity analyses focus on the biological distance or
35
36 similarity between the interacting proteins. Such distance/similarity based matrix can then be
37
38 utilised to build up the traditional clustering algorithms as in, for instance, Unweighted Pair
39
40 Group Method with Arithmetic Mean (UPGMA), generally used for calculating evolutionary
41
42 distance. However, to emerge into more biologically relevant models, instead of only indicating
43
44 the binary relationships as in the traditional coefficient based ones, sequence similarity, structural
45
46 similarity and gene expression correlation have started to be used [55-57].
47
48
49
50
51

52 Any attempt to cluster such biologically relevant modular networks would bring out the
53
54 importance of the interrelationships of the constituent components. To formulate the modularity,
55
56
57
58
59
60

1
2
3 the graph theoretic and the topology-based methods consider the local or global structure of the
4
5 PPI networks. While the former converts the process of clustering into graph theoretic problems,
6
7
8 the latter quantitatively measures the metric features of the networks before formulating the
9
10 clustering algorithms for modularity analyses. It is to be noted that the graph theoretic features
11
12 have gained much importance in modularity analyses due to the fact that they can find out the
13
14 densest subnetworks e.g. Molecular Complex Detection (MCODE), clique percolation. Amongst
15
16 these, clique percolation method has its advantage of identifying overlapping functional clusters
17
18 in a typical PPI network. This enables one to detect proteins simultaneously functioning
19
20 differently in several different modules [3]. The other method of utilising the graph theoretic
21
22 measures is through partitioning the modular subnetworks, either by simple partition detection
23
24 through less important edges or by an improved Markov clustering algorithm which uses the
25
26 mathematical bootstrapping procedure [3].
27
28
29
30

31 One of the recent methods entails a flow-based technique which can deal both with the
32
33 prediction of protein function and protein modularity analysis. There are several algorithms
34
35 which have been developed with this concept. One of them is the ‘Majority’ method which
36
37 considers the interactions of its neighbors and adopts the three most frequent annotations [58].
38
39 An extension of the above method, ‘Neighborhood’, employs a search for all the proteins within
40
41 a particular radius to identify overrepresented functional annotations [59]. The usage of edge
42
43 weights through gene expression data was done by Karaoz *et al.* [60]. Similar kind of weighted
44
45 interaction network was used following a ‘guilt-by-association’ principle, wherein the functional
46
47 flow was created from the annotated protein to the unannotated ones, through simulation [61].
48
49 Such kind of simulation of biological or functional flows within the network can be used as an
50
51 essential tool of modeling to explore the dynamic signal transduction systems [3]. Moreover,
52
53
54
55
56
57
58
59
60

1
2
3 network flow simulations can predict complex network behavior under a realistic variety of
4 external stimuli. A very important algorithm called CASCADE helps to detect the dynamic flow
5 simulation of modularity analyses. CASCADE utilises the concept of occurrence probability and
6 models a unique clustering methodology encompassing the biological and topological influence
7 of each protein on the other. Occurrence probability brings out the distribution of the number of
8 interactions necessary to link a pair of instant proteins in the network at a given time point [62].
9

10 The methods for the generation and analyses of the networks discussed as of now would be more
11 accurate with a benchmarking of the data. Clustering techniques described here are based solely
12 upon the graph theoretical properties without any real supervised data, **thereby confirming** their
13 authenticity. However, *a priori* knowledge from amino acid and genomic sequences, protein
14 structures and evolutionary profiles, gene expression and ontology annotation could be integrated
15 with the PPI data to add to the analyses. Information about protein domains and localization has
16 been used to successfully predict protein functions [63, 64]. A variety of high throughput data
17 including microarray and protein complex data have been integrated to construct Bayesian
18 models [65, 66], and Kernel based matrices have also been proposed [67, 68].
19

20 It is worth mentioning at this point that different clustering techniques and even the same
21 technique with different parameters end up in giving disparate outcomes. Thus, validation of
22 these clustering techniques is mandatory. Indeed, different clustering algorithms have been
23 evaluated by several researchers in order to understand their potential to infer protein clusters
24 from protein interaction networks [69, 70]. Jiang *et al.* [71], Zhang [3] have suggested different
25 approaches to validate clustering methods, including validation based on agreement with
26 annotated protein function databases, definition of clustering, the reliability of clusters,
27 topological properties and the p-value from the hypergeometric distribution.
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

The visualization

There are several software and plug-ins, added therein, for the visualization of the real networks constructed in the form of graphs of the interconnected proteins. Researchers across the world have used Cytoscape (latest version 2.8.2) [72] and Gephi (latest version beta 0.8.2) [73]. Cytoscape has the plug-in, NETWORK ANALYZER [74], to compute values for the classical network centrality parameters like DC, CC and BC besides clustering coefficient, average path length (APL) and network diameter. Another important centrality measure, the EC, can be calculated via Gephi. The Java plug-in, cytoHubba [75], can be used to categorise the top ranked proteins/hubs in the network. Combined scores, from different parameters considered in the databases like STRING, can be taken as edge weights for computing Cytohubba scores. Several topological algorithms, viz. Maximal Clique Centrality (MCC), Maximum Neighborhood Component (MNC), and Density of Maximum Neighborhood Component (DMNC), can be used to find the important hub proteins of the networks. To obtain the clusters of proteins after the network decomposition, the Molecular Complex Detection (MCODE) algorithm can be implemented to find the densely connected regions in the networks [52].

Analyses of pathogenic PINs: Intra- and Interspecies scenario

In order to gain insight into the infection strategies of pathogens, several intra-pathogenic and host-pathogen protein interaction networks have been generated and analysed over the last decade. This section will delineate the scenario of protein interaction network analysis of some of these species including viruses, bacteria and protozoan parasites (Table 1). Amongst these, the topology of intra-viral networks of different members of herpesvirus family (*viz.*, Kaposi's

1
2
3 sarcoma-associated herpesvirus (KSHV), Varicella-zoster virus (VZV), Epstein-Barr virus
4 (EBV)), SARS-coronavirus (SARS-CoV), Hepatitis C virus (HCV) and Influenza A virus (H1N1
5 and H3N2) have been investigated by evaluating different network parameters like degree, APL,
6 clustering coefficient and network diameter [76-82]. The analysis revealed that viral networks
7 appear as single, highly coupled modules with relatively many hubs and few 'peripheral' nodes,
8 in contrast to scale-free cellular networks having well-separated functional modules. This
9 distinguishing network topology, may be essential for the formation of compact virions and
10 functional viral complexes. However, it is unclear whether the disparity between viral and
11 cellular network topology is a consequence of biological differences or the artifacts of
12 experimental biases and errors.
13
14
15
16
17
18
19
20
21
22
23
24
25

26
27 Furthermore, the comparison of interactomes can lead to the identification of highly conserved
28 interactions, critical for pathogenesis and thus, could serve as promising broad spectrum drug
29 targets. For example, the comparison of intra-viral networks for herpesviruses enabled to identify
30 a core set of highly conserved interactions, which mediate budding of capsids at the inner nuclear
31 membrane of the host, and thus, could be promising targets for alternative herpesviral therapies.
32
33
34
35
36
37

38
39 The first large-scale intra-bacterial PPI networks were constructed and analyzed for *Helicobacter*
40 *pylori* and subsequently, for several other bacterial pathogens, such as *Campylobacter jejuni*,
41 *Treponema pallidum*, *Mycoplasma pneumonia*, *Mycobacterium tuberculosis*, and *Staphylococcus*
42 *aureus* [44, 83-87]. The topological parameters (degree and BC) of undirected intra-bacterial
43 networks, studied so far, revealed that bacterial networks are scale-free in nature, following a
44 power law distribution. The evaluation of average clustering coefficient of bacterial protein
45 networks (eg., *C.jejuni*, *M. tuberculosis* etc.) indicated that networks comprised of many clusters
46 ie., subnetworks of highly interconnected proteins and comparative network analysis (CNA)
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 revealed that many of the subnetworks were conserved across organisms, identified using the
4 NetworkBlast algorithm [84, 87]. It is worth to be noted that the proteins enriched in conserved
5 subnetworks carry out specific Gene Ontology (GO) functions representing crucial functional
6 pathways or protein complexes. Indeed, in *C. jejuni*, clustering of the conserved subnetworks
7 using k-means algorithm followed by UPGMA identified core proteins having distinct cellular
8 function. These core proteins were found to present in many subnetworks [84]. Thus, the
9 organism's interaction network can be used to predict the function of the unannotated proteins or
10 to map protein complexes and pathways involved in virulence, providing the directions for
11 uncovering new drug targets [44, 84, 87]. Network topology was exploited to identify essential
12 genes/proteins, which are crucial for replication, growth and viability of an organism, in different
13 pathogenic species, including *S. aureus*, *C.jejuni*, *M. tuberculosis*, *Mycobacterium abscessus*,
14 and various food and waterborne pathogens [44, 84, 88-90]. Proteins, encoded by essential
15 genes, are hub proteins with many number of interactions in a network, and are also important
16 for network integrity and stability, thus could be potential to be therapeutic targets.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

36 The protein interaction networks for the above mentioned viruses, including *Human*
37 *Immunodeficiency Virus* (HIV)-1 [76, 77, 81, 91-93] and different bacterial pathogens, such as
38 *Bacillus anthracis*, *Francisella tularensis*, and *Yersinia pestis*, and *M.tuberculosis* with their
39 human host have been studied [94-96]. The network topology analyses of host-pathogen systems
40 indicated that both viral and bacterial proteins target human proteins which own higher degree
41 and higher BC in the human protein interaction network. Viruses and bacteria both follow a
42 common infection strategy of preferentially attacking hub and bottleneck proteins to impede
43 host's essential biology [82]. Viruses tend to interact with host proteins which have higher
44 degree and BC values compared to their bacterial counterparts. Identification of conserved
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 subnetworks in human-pathogen PPI [94] and GO functional analysis of pathogen-targeted
4 human proteins, delineated a perfect picture of their infection strategy. Bacteria upholds infection
5 in humans by foraying proteins involved in immune response thereby shattering human defense
6 mechanism, whereas viruses exploit host's transcriptional machinery to propagate themselves
7 within the host. It is worthy to mention at this point that most of the pathogen-targeted host
8 proteins are those that play critical role in regulation of metabolic processes, such as cell-cycle
9 regulation, nuclear transport and most importantly immune response.

10
11 *Plasmodium falciparum*, the causative agent of malaria in human, is the only protozoan parasite
12 whose protein interaction network has been studied extensively [97-100]. Each study mostly
13 focused on the identification and isolation of critical protein clusters/subnetworks or pathways,
14 and also assigning the function of uncharacterized proteins. The study identified a group of
15 proteins, such as chaperones, transcription factors and new surface proteins which are crucial for
16 parasite's invasion and survival. Most of the proteins in the highly interconnected subnetworks
17 were found to be involved in pathogenesis, perhaps the result of gene duplication event for
18 maintaining its parasitic way of life. The identification of subnetworks was mainly done by using
19 clustering coefficient, Markov clustering algorithm [97], clique percolation algorithm [98] and k-
20 means clustering. It is worth mentioning at this point that plasmodium network stands distinct
21 from other eukaryotic network because of its 'assortative' nature and bearing very less overlap
22 with their interactomes. A very recent study [100] aimed at identifying important interacting
23 proteins (IIPs) in Plasmodium PPI network, using various node centrality indices (degree,
24 closeness, radiality, betweenness, eccentricity, stress, weinner index, centroid, assortativity and
25 clustering coefficient) and network centrality indices (average distance, connectivity distribution,
26 diameter and average clustering coefficient), followed by *in silico* knock-out analysis. The
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 highly interacting hub and central proteins, which are vital for network integrity as well as
4
5 crucial for organism's survival, were considered as important proteins. These IIPs also play a
6
7 vital role in stage specificity and were found to interact with several human proteins associated
8
9 with multiple metabolic pathways, signaling pathways and infection mechanism. It has to be
10
11 noted that human proteins targeted by pathogen are hubs in the human interactome and
12
13 malfunctioning of the crucial host pathways results in clinical manifestation of malaria, which
14
15 pose the interacting pathogen proteins as potential drug targets.
16
17
18
19
20
21

22 **Applications in systems biomedicine**

23
24
25 With the main target of the present review being the application of the PPI networks in
26
27 biomedicine, cases to unravel the molecular basis of disease, by studying disease related
28
29 subnetworks, have been reported. For instance, a new dimension has been given by such PPI
30
31 network analysis to bring out the relationship of the pathogenic bacterium *Helicobacter pylori*
32
33 with gastric carcinoma [101]. This has achieved a level of acceptance from the World Health
34
35 Organization (WHO) and the International Agency for Research on Cancer consensus groups
36
37 who have classified *H. pylori* as a definite biological carcinogen. The authors analysed the
38
39 networks built upon the selected translated proteins of the expressed genes from databases and
40
41 literatures. Their analyses reflected connectors of oncogenic proteins as hub and bottleneck
42
43 proteins, mostly related to immune response governing the cell cycle, cell maintenance and
44
45 proliferation, and transcription regulators [101].
46
47
48
49
50

51 An indirect study on a smaller scale upon *Salmonella* Pathogenicity Island Type III secretion
52
53 system was carried out to build a methodology of targeting the indispensable proteins from
54
55 amongst a conglomerate [42]. The authors constructed the network from the available
56
57
58
59
60

1
2
3 interactions from STRING database and analysed it with the common and rarely used centrality
4
5 measures to decide upon the most indispensable one. They benchmarked their theoretical finding
6
7 through analyses of networks built from the expressed gene products of two different microarray
8
9 data and arrived at the same point with respect to such indispensable protein issue [42]. The
10
11 outcome of these two works clearly would be the positive side of the analyses of PPI networks
12
13 for generating systems biomedicine where the goal would be to harm the pathogen without
14
15 harming the host and avoiding rapid development of antimicrobial resistance. The discussion on
16
17 such issues takes us to a point wherein workers in this field would like to keep in mind few
18
19 points while carrying out the related research. As indicated earlier, Lahiri *et al.* [42] have
20
21 delineated a key methodology which could be followed with modifications as and when needed.
22
23 A network constructed from a source has to be checked in for scale-freeness. The network can
24
25 then be pruned to a core of proteins and/or top rankers from different centrality measures can be
26
27 compared to unanimity. The finding therefrom can be benchmarked by other experimental omics
28
29 data to corroborate. A selection of centrality measures would depend upon the requirement of the
30
31 work. Following just some network analyses and trying to get a positive outcome, however,
32
33 would abrogate the essentiality of PIN analyses.
34
35
36
37
38
39
40
41
42

43 **Dependability of analyzed PINs**

44
45 While there can be claims about the necessity of PIN analyses, a very important point needs to be
46
47 considered to facilitate such claims. It is to be understood that the correctness of the analyses of
48
49 such PIN would depend upon the correct construction of the network. Many such networks are
50
51 being built based upon laboratory experimentation like yeast two-hybrid and mass spectrometry
52
53 data generation. Moreover, networks built from various sources have extremely low overlap of
54
55
56
57
58
59
60

1
2
3 different high-throughput data generating manually curated databases. The other possibilities
4 causing such error, down the line of network analyses, could be low reliability of literature
5 curation and difficulties arising due to improper gene annotation and webpage data extraction
6 [102]. As the above methods can be highly error prone, the dependability of PIN analyses
7 become low [103]. In fact, there can be falsely reported interactions as well as left out
8 interactions not being reported. Alarming false discovery rates (FDRs) of 10-20% and false
9 negative rates (FNRs) of upto 50% are reported for yeast, worm and fly screens [104, 105].
10 However, such falsification of interactions could also crop from the low coverage of different
11 comparative methods having noises leading to misinterpretation and **erroneous** integration of
12 data [106]. An interesting concept on such comparative methods of interactions is (~~to note~~) that a
13 comparison of the individual proteins interaction reveals a common tendency between methods
14 manifested as global properties of the PINs [107]. To reduce such uncertainties of PIN
15 construction from experimental data, two models have been proposed. These are the spoke and
16 the matrix model of studying the bait connecting the prey. The former, connecting the bait along
17 with the hit proteins, yields less false positives and is three times more accurate than the matrix
18 model which connects all proteins. However, the latter yields more true positives as well [108].
19 A list of such sources of PIN can be **obtained** from literatures [102, 109].

20 Irrespective of the network construction, the analyses, however, can be of potential in cases of
21 assessing the efficacy of a drug target, where a specific pathway is targeted to inhibit it. In this
22 case, a perturbation of a dynamic network by inhibiting a specific pathway is manifested as a
23 diversion to alternate pathways, as discussed in CASCADE [3]. However, the shortest path
24 distance between important proteins **remains** the same and thus, proteins connecting **other**
25 important **ones** in the network and thereby bridging them, have high BC. Instances of reduction
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 of such alternate pathways, keeping the core pathway intact as the shortest path, in metabolic
4 network of *Mycobacterium leprae*, have been reported [110].
5
6
7
8
9

10 **Concluding remarks**

11
12 The present review has delineated a broad overview of analysing protein interaction networks of
13 infectious diseases caused by viruses, bacteria and protozoan parasites. It entails the different
14 methodologies which can be adopted by researchers while trying to analyse such networks. A
15 thorough look of the review shows that most of the researchers resorted to only a handful of the
16 techniques to conclude about important protein identification, pathway detection and functional
17 prediction. It is imperative, however, that a benchmarking of these computationally predicted
18 and analysed results would be mandatory for a better future towards non-conventional health
19 intervention processes. For instance, Lahiri *et al.* [42] adopted several of these techniques and
20 then validated with some biologically relevant high throughput microarray data. Days are not far
21 when it would be a practice for the researchers to spread themselves and come up with new
22 health intervention strategies to generate more accurate systems-based biomedicines.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

41 **Acknowledgement**

42 Authors AP, AR and BS are thankful to Centre for Bioinformatics, Pondicherry University,
43 Pondicherry, India and CL to the Department of Biological Sciences, Sunway University,
44 Selangor, Malaysia for proving infrastructure facilities to carry on the work. AR and BS are
45 indebted to Pondicherry University for their pre-doctoral fellowship.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

1. Aderem A, Adkins JN, Ansong C, *et al.* A systems biology approach to infectious disease research: innovating the pathogen-host research paradigm. *MBio* 2011;**2**:e00325-10.
2. Maragakis LL, Perl TM. *Acinetobacter baumannii*: epidemiology, antimicrobial resistance, and treatment options. *Clin Infect Dis* 2008;**46**:1254-63.
3. Zhang A. Protein Interaction Networks: Computational Analysis. New York: Cambridge University Press, 2009.
4. Uetz P, Giot L, Cagney G, *et al.* A comprehensive analysis of protein-protein interactions in *Saccharomyces cerevisiae*. *Nature* 2000;**403**:623-7.
5. Ho Y, Gruhler A, Heilbut A, *et al.* Systematic identification of protein complexes in *Saccharomyces cerevisiae* by mass spectrometry. *Nature* 2002;**415**:180-3.
6. Ge H., UPA, a universal protein array system for quantitative detection of protein-protein, protein-DNA, protein-RNA and protein-ligand interactions. *Nucleic Acids Res* 2000;**28**:e3.
7. Marcotte EM, Pellegrini M, Ng HL, *et al.* Detecting protein function and protein-protein interactions from genome sequences. *Science* 1999;**285**:751-3.
8. Overbeek R, Fonstein M, D'Souza M, *et al.* The use of gene clusters to infer functional coupling. *Proc Natl Acad Sci USA* 1999;**96**:2896-901.
9. Matthews LR, Vaglio P, Reboul J, *et al.* Identification of potential interaction networks using sequence-based searches for conserved protein-protein interactions or "Interologs". *Genome Res* 2001;**11**:2120-6.
10. Wojcik J, Schächter V. Protein-protein interaction map inference using interacting domain profile pairs. *Bioinformatics* 2001;**17**:S296-S305.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
11. Aloy P, Russell RB. Interrogating protein interaction networks through structural biology. *Proc Natl Acad Sci USA* 2002;**99**:5896-901.
12. Smith GR, Sternberg MJ. Prediction of protein-protein interactions by docking methods. *Curr Opin Struct Biol* 2002;**12**:28-35.
13. Bock JR, Gough DA. Predicting protein-protein interactions from primary structure. *Bioinformatics* 2001;**17**:455-60.
14. Jansen R, Yu H, Greenbaum D, *et al.* A Bayesian networks approach for predicting protein-protein interactions from genomic data. *Science* 2003;**302**:449-53.
15. Prasad TSK, Goel R, Kandasamy K, *et al.* Human protein reference database- update 2009. *Nucleic Acids Res* 2009;**37**:D767-72.
16. Persico M, Ceol A, Gavrila C, *et al.* HomoMINT: an inferred human network based on orthology mapping of protein interactions discovered in model organisms. *BMC Bioinformatics* 2005;**6**:S21.
17. Licata L, Briganti L, Peluso D, *et al.* MINT, the molecular interaction database: 2012 update. *Nucleic Acids Res* 2012;**40**:D857-61.
18. Stark C, Breitkreutz BJ, Reguly T, *et al.* BioGRID: a general repository for interaction datasets. *Nucleic Acids Res* 2006;**34**:D535-9.
19. Salwinski L, Miller CS, Smith AJ, *et al.* The database of interacting proteins: 2004 update. *Nucleic Acids Res* 2004;**32**:D449-51.
20. Pagel P, Kovac S, Oesterheld M, *et al.* The MIPS mammalian protein-protein interaction database. *Bioinformatics* 2005;**21**:832-4.
21. Alfarano C, Andrade CE, Anthony K, *et al.* The Biomolecular Interaction Network Database and related tools 2005 update. *Nucleic Acids Res* 2005;**33**:D418-24.

- 1
2
3 22. Franceschini A, Szklarczyk D, Frankild S, *et al.* STRING v9.1: protein-protein
4 interaction networks, with increased coverage and integration. *Nucleic Acids Res* 2013;
5
6 41:D808-15.
7
8
- 9
10 23. Kerrien S, Aranda B, Breuza L, *et al.* The IntAct molecular interaction database in 2012.
11
12 *Nucleic Acids Res* 2012;40:D841-6.
13
14
- 15 24. Warde-Farley D, Donaldson SL, Comes O, *et al.* The GeneMANIA prediction server:
16 biological network integration for gene prioritization and predicting gene function.
17
18 *Nucleic Acids Res* 2010;38:W214-20.
19
20
- 21 25. Schmitt T, Ogris C, Sonnhammer EL. FunCoup 3.0: database of genome-wide functional
22 coupling networks. *Nucleic Acids Res* 2014;42:D380-8.
23
24
25
26
- 27 26. Kamburov A, Stelzl U, Lehrach H, *et al.* The ConsensusPathDB interaction database:
28 2013 update. *Nucleic Acids Res* 2013;41:D793-800.
29
30
- 31 27. Wattam AR, Abraham D, Dalay O, *et al.* PATRIC, the bacterial bioinformatics database
32 and analysis resource. *Nucleic Acids Res* 2014;42:D581-91.
33
34
35
- 36 28. Rid R, Strasser W, Siegl D, *et al.* PRIMOS: an integrated database of reassessed protein-
37 protein interactions providing web-based access to in silico validation of experimentally
38 derived data. *Assay Drug Dev Technol* 2013;11:333-46.
39
40
41
42
- 43 29. Kumar R, Nanduri B. HPIDB-a unified resource for host-pathogen interactions. *BMC*
44 *Bioinformatics* 2010;11:S16.
45
46
47
- 48 30. Urban M, Pant R, Raghunath A, *et al.* The Pathogen-Host Interactions database (PHI-
49 base): additions and future developments. *Nucleic Acids Res* 2015;43:D645-55.
50
51
52
- 53 31. Winnenburger R, Baldwin TK, Urban M, *et al.* PHI-base: a new database for pathogen host
54 interactions. *Nucleic Acids Res* 2006;34:D459-64.
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
32. Calderone A, Licata L, Cesareni G. VirusMentha: a new resource for virus-host protein interactions. *Nucleic Acids Res* 2015;**43**:D588-92.
 33. Guirimand T, Delmotte S, Navratil V. VirHostNet 2.0: surfing on the web of virus/host molecular interactions data. *Nucleic Acids Res* 2015;**43**:D583-7.
 34. Erdos P, Renyi A. On the evolution of random graphs, Publications of Mathematical Institute of Hungarian Academy of Science, 1960;**5**:17-61.
 35. Watts DJ, Strogatz SH. Collective dynamics of 'small-world' networks. *Nature* 1998;**393**:440-2.
 36. Albert R, Jeong H, Barabasi AL. Error and attack tolerance of complex networks. *Nature* 2000;**406**:378-482.
 37. Jeong H, Mason SP, Barabasi AL, *et al.* Lethality and centrality in protein networks. *Nature* 2001;**411**:41-2.
 38. Mason O, Verwoerd M. Graph theory and networks in Biology. *IET Syst Biol* 2007;**1**:89-119.
 39. Pavlopoulos GA, Hooper SD, Sifrim A, *et al.* A tool for exploring and clustering biological networks. *BMC Res Notes* 2011;**4**:384.
 40. Ozgur A, Vu T, Erkan G, *et al.* Identifying gene-disease associations using centrality on a literature mined gene-interaction network. *Bioinformatics* 2008;**24**:277-85.
 41. Koschützki D, Schreiber F. Comparison of centralities for biological networks. In: Proceedings of the German Conference on Bioinformatics (GCB), 2004;**53**:199-206.
 42. Lahiri C, Pawar S, Sabarinathan R, *et al.* Interactome analyses of Salmonella pathogenicity islands reveal SicA indispensable for virulence. *J Theor Biol* 2014;**363**:188-97.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
43. Bonacich P. Factoring and weighting approaches to status scores and clique identification. *Journal of Mathematical Sociology* 1972;**2**:113-120.
 44. Cherkasov A, Hsing M, Zoraghi R, et al. Mapping the protein interaction network in Methicillin-resistant *Staphylococcus aureus*, *J Proteome Res* 2011;**10**:1139-1150.
 45. Dickerson J, Pinney J, Robertson D. The biological context of HIV-1 host interactions reveals subtle insights into a system hijack. *BMC Syst Biol* 2010;**4**:1752.
 46. Hahn MW, Kern, AD. Comparative genomics of centrality and essentiality in three eukaryotic protein-interaction networks. *Mol Biol Evol* 2005;**22**:803-6.
 47. Han J, Bertin N, Hao T, et al. Evidence for dynamically organized modularity in the yeast protein-protein interaction network. *Nature* 2004;**430**:88-93.
 48. Hwang W, Cho YR, Zhang A, et al. A novel functional module detection algorithm for protein-protein interaction networks. *Algorithms Mol Biol* 2006;**1**:24.
 49. Dartnell L, Simeonidis E, Hubank M, et al. Robustness of the p53 network and biological hackers. *FEBS Lett* 2005;**579**:3037-42.
 50. Xu K, Bezakova I, Bunimovich L, et al. Path lengths in protein-protein interaction networks and biological complexity. *Proteomics* 2011;**11**:1857-67.
 51. Asif W, Qureshi A, Iqbal M, et al. On the complexity of average path length for biological networks and patterns. *Int J Biomath* 2014; 7.
 52. Bader GD, Hogue CW. An automated method for finding molecular complexes in large protein interaction networks. *BMC Bioinformatics* 2003;**4**:2.
 53. Seidman SB. Network structure and minimum degree. *Social Networks* 1983; 5:269-87.
 54. Bollobas B. The evolution of sparse graphs. In *Graph Theory and Combinatorics*, Proceeding Cambridge Combinatorial Conference in honor of Paul Erdos, 1984, 35-57.

- 1
2
3 55. Enright AJ, van Dongen S, Ouzounis CA. An efficient algorithm for large-scale detection
4 of protein families. *Nucleic Acids Res* 2002;**30**:1575-84.
5
6
- 7
8 56. Domingues F, Rahnenfuhrer J, Lengauer T. Automated clustering of ensembles of
9 alternative models in protein structure databases. *Protein Eng Des Sel* 2004;**17**:537-43.
10
11
- 12 57. Tornow S, Mewes HW. Functional modules by relating protein interaction networks and
13 gene expression. *Nucleic Acids Res* 2003;**31**:6283-9.
14
15
- 16 58. Schwikowski B, Uetz P, Fields S. A network of protein-protein interactions in yeast. *Nat*
17 *Biotechnol* 2000;**18**:1257-61.
18
19
- 20 59. Hishigaki H, Nakai K, Ono T, *et al.* Assessment of prediction accuracy of protein
21 function from protein-protein interaction data. *Yeast* 2001;**18**:523-31.
22
23
- 24 60. Karaoz U, Murali TM, Letovsky S, *et al.* Whole-genome annotation by using evidence
25 integration in functional-linkage networks. *Proc Natl Acad Sci USA* 2004;**101**:2888-93.
26
27
- 28 61. Nabieva E, Jim K, Agarwal A. Whole-proteome prediction of protein function via graph-
29 theoretic analysis of interaction maps. *Bioinformatics* 2005;**21**:i302-10.
30
31
- 32 62. Hwang W, Kim T, Ramanathan M, *et al.* Bridging centrality: graph mining from element
33 level to group level. In Proceedings of the 14th ACM SIGKDD International Conference
34 on Knowledge Discovery & Data Mining (KDD08), 2008;336-44.
35
36
- 37 63. Chen X, Liu M, Ward R. Protein function assignment through mining cross-species
38 protein-protein interactions. *PLoS One* 2008;**3**:e1562.
39
40
- 41 64. Nariai N, Kasif S. Context specific protein function prediction. *Genome Inform* 2007;
42 **18**:173-82.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
65. Troyanskaya OG, Dolinski K, Owen AB, *et al.* A Bayesian framework for combining heterogeneous data sources for gene function prediction. *Proc Natl Acad Sci USA* 2003; **100**:8348-53.
 66. Chen Y, Xu D. Global protein function annotation through mining genome-scale data in yeast *Saccharomyces cerevisiae*. *Nucleic Acids Res* 2004;**32**:6414-24.
 67. Lanckriet GRG, Deng M, Cristianini N, *et al.* Kernel-based data fusion and its application to protein function prediction in yeast. Pacific Symposium on Biocomputing, 2004;9.
 68. Tsuda K, Shin HJ, Schoelkopf B. Fast protein classification with multiple networks. 2005;**21**:ii59-65.
 69. Brohee S, van Helden J. Evaluation of clustering algorithms for protein-protein interaction networks. *BMC Bioinformatics* 2006;**7**:488.
 70. Moschopoulos CN, Pavlopoulos GA, Iacucci E, *et al.* Which clustering algorithm is better for predicting protein complexes? *BMC Res Notes* 2011;**4**:549.
 71. Jiang D, Tang C, Zhang A. Cluster analysis for gene expression data: A Survey. *IEEE Transactions on Knowledge and Data Engineering (TKDE)* 2004;**16**:1370-86.
 72. Shannon P, Markiel A, Ozier O, *et al.* Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res* 2003;**13**:2498-504.
 73. Bastian M, Heymann S, Jacomy M. Gephi: an open source software for exploring and manipulating networks. In: Proceedings of the International AAAI Conference on Weblogs and Social Media. San Jose, CA, North America (ICWSM09), 2009.
 74. Assenov Y, Ramirez F, Schelhorn SE, *et al.* Computing topological parameters of biological networks. *Bioinformatics* 2008;**24**:282-4.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
75. Lin CY, Chin CH, Wu HH, *et al.* Hubba: hub objects analyzer-a framework of interactome hubs identification for network biology. *Nucleic Acids Res* 2008;**36**:W438-43.
76. Uetz P, Dong YA, Zeretzke C, *et al.* Herpesviral protein networks and their interaction with the human proteome. *Science* 2005;**311**:239-42.
77. Calderwood MA, Venkatesan K, Xing L, *et al.* Epstein-Barr virus and virus human protein interaction maps. *Proc Natl Acad Sci USA* 2007;**104**:7606-11.
78. Bailer SM, Hass J. Connecting viral with cellular interactomes. *Curr Opin Microbiol* 2009;**12**:453-9.
79. von Brunn A, Teepe C, Simpson JC, *et al.* Analysis of intraviral protein-protein interactions of the SARS coronavirus ORFeome. *PLoS One* 2007;**2**:e459.
80. Meyniel-Schicklin L, de Chasseay B, André P, *et al.* Viruses in interactomes in translation. *Mol Cell Proteomics* 2012;11.
81. Shapira SD, Gat-Viks I, Shum BO, *et al.* A physical and regulatory map of host-influenza interactions reveals pathways in H1N1 infection. *Cell* 2009;**139**:1255-67.
82. Tekir SD, Çakır T and Ülgen KÖ. Infection strategies of bacterial and viral pathogens through pathogen-human protein-protein interactions. *Front in Microbiol*, 2012, **3**:46,1-11.
83. Rain JC, Selig L, De Reus H, *et al.* The protein-protein interaction map of *Helicobacter pylori*. *Nature* 2001;**409**:211-5.
84. Parrish JR, Yu J, Liu G, *et al.* A proteome-wide protein interaction map for *Campylobacter jejuni*. *Genome Biol* 2007;**8**:R130.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
85. Titz B, Rajagopala SV, Goll J, *et al.* The binary protein interactome of *Treponema pallidum* - the syphilis spirochete. *PLoS One* 2008;**3**:e2292.
86. Kühner S, van Noort V, Betts MJ, *et al.* Proteome organization in a genome-reduced bacterium. *Science* 2009;**326**:1235-40.
87. Wang Y, Cui T, Zhang C, *et al.* Global protein-protein interaction network in the human pathogen *Mycobacterium tuberculosis* H37Rv. *J Proteome Res* 2010;**9**:6665-77.
88. Raman K, Yeturu K, Chandra N. targetTB: a target identification pipeline for *Mycobacterium tuberculosis* through an interactome, reactome and genome-scale structural analysis. *BMC Syst Biol* 2008;**2**:109.
89. Shanmugham B, Pan A. Identification and characterization of potential therapeutic candidates in emerging human pathogen *Mycobacterium abscessus*: a novel hierarchical in silico approach. *PLoS One* 2013;**8**:e59126.
90. Jadhav A, Shanmugham B, Rajendiran A, *et al.* Unraveling novel broad-spectrum antibacterial targets in food and waterborne pathogens using comparative genomics and protein interaction network analysis. *Infect Genet Evol* 2014;**27**:300-8.
91. de Chasse B, Navratil V, Tafforeau L, *et al.* Hepatitis C virus infection protein network. *Mol Syst Biol* 2008;**4**:230.
92. Emamjomeh A, Goliaei B, Zahirab J, *et al.* Predicting protein-protein interactions between human and hepatitis C virus via an ensemble learning method. *Mol Biosyst* 2014;**10**:3147-54.
93. Bandyopadhyay S, Ray S, Mukhopadhyay A, *et al.* A review of in silico approaches for analysis and prediction of HIV-1-human protein-protein interactions. *Brief Bioinform* 2014.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
94. Dyer MD, Neff C, Dufford M, *et al.* The human-bacterial pathogen protein interaction network of *Bacillus anthracis*, *Francisella tularensis*, and *Yersinia pestis*. *PLoS One* 2010; **5**:e12089.
95. Yang H, Ke Y, Wang J, *et al.* Insight into bacterial virulence mechanisms against host immune response via the *Yersinia pestis*-human protein-protein interaction network. *Infect Immun* 2011;**79**:4413-24.
96. Zhou H, Gao S, Nguyen NN, *et al.* Stringent homology-based prediction of *H. sapiens*-*M. tuberculosis* H37Rv protein-protein interactions. *Biol Direct* 2014;**9**:5.
97. Wuchty S, Ipsaro JJ. A Draft of Protein Interactions in the Malaria Parasite *P. falciparum*. *J Proteome Res* 2007;**6**:1461-70.
98. Wuchty S, Adams JH, Ferdig MT. A comprehensive *Plasmodium falciparum* protein interaction map reveals a distinct architecture of a core interactome. *Proteomics* 2009;**9**: 1841-9.
99. Wuchty S. Computational Prediction of Host-Parasite Protein Interactions between *P. falciparum* and *H. sapiens* *PLoS One* 2011;**6**:e26960.
100. Bhattacharyya M, Chakrabarti S. Identification of important interacting proteins (IIPs) in *Plasmodium falciparum* using large-scale interaction network analysis and in silico knock-out studies. *Malaria Journal* 2015;**14**:70.
101. Kim KK, Kim HB. Protein interaction network related to *Helicobacter pylori* infection response. *World J Gastroenterol* 2009;**15**:4518-28.
102. Cusick ME, Yu H, Smolyar A, *et al.* Literature-curated protein interaction datasets. *Nat Methods* 2009;**6**:39-46.

- 1
2
3 103. Rinner O, Mueller LN, Hubálek M, *et al.* An integrated mass spectrometric and
4 computational framework for the analysis of protein interaction networks. *Nat Biotechnol*
5 2007;**25**:345-52.
6
7
8
9
10 104. Huang H, Bader JS. Precision and recall estimates for two-hybrid screens.
11 *Bioinformatics* 2009;**25**:372-8.
12
13
14
15 105. Hart GT, Ramani AK, Marcotte EM. How complete are current yeast and human protein-
16 interaction networks? *Genome Biol* 2006;**7**:120.
17
18
19
20 106. Gentleman R, Huber W. Making the most of high-throughput protein-interaction data.
21 *Genome Biol* 2007;**8**:112.
22
23
24
25 107. Hoffmann R, Valencia A. Protein interaction: same network, different hubs. *Trends*
26 *Genet* 2003;**19**:681-3.
27
28
29
30 108. Bader GD, Hogue CW. Analyzing yeast protein-protein interaction data obtained from
31 different sources. *Nat Biotechnol* 2002;**20**:991-7.
32
33
34
35 109. Sanderson CM. The cartographers toolbox: building bigger and better human protein
36 interaction networks. *Brief Funct Genomic and Proteomic* 2009;**8**:1-11.
37
38
39 110. Verkhedkar KD, Raman K, Chandra NR, *et al.* Metabolome Based Reaction Graphs of
40 M. tuberculosis and M. leprae: A Comparative Network Analysis. *PLoS One*
41 2007;**2**:e881.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table legend

Table1: Mode of protein interaction network analysis in intra-pathogenic and host-pathogen systems

Supplementary Table1: The number of intra-pathogenic and host-pathogen PPIs from different resources along with pathogens' genome size

Figure legend

Figure 1: Schematic representation summarizing different computational approaches to analyze protein interaction networks.

Table 1: Mode of protein interaction network analysis in intra-pathogenic and host-pathogen systems

Pathogen	Disease	Mode of Network analysis	References
Kaposi's sarcoma-associated herpesvirus (KSHV) ^{φ,ψ}	Kaposi sarcoma, B-cell lymphomas	^φ DC,BC,DD,APL,CCf,ND ^ψ DC	[76,80]
Varicella zoster virus (VZV) ^{φ, ψ}	Chickenpox, shingles	^φ DC,BC,DD,APL,CCf,ND ^ψ DC	[76,80]
Epstein-Barr virus (EBV) ^{φ,ψ}	Mononucleosis	^φ DC,BC,DD,APL,CCf,ND ^ψ DC,APL,CCf	[77,80]
Severe acute respiratory syndrome-coronavirus (SARS-CoV) ^φ	Severe acute respiratory syndrome	^φ DC,BC,DD,APL,CCf,ND	[79,80]
Hepatitis C virus (HCV) ^{φ,ψ}	Hepatitis	^φ DC,BC,DD,APL,CCf,ND ^ψ DD,BC,CC,APL,CCf,NC,Str,Ecc,Ra,PA(KEGG),GO,DA	[80,91,92]
Influenza A virus (H1N1, H3N2) ^{φ,ψ}	Influenza	^φ DC,BC,DD,APL,CCf,ND ^ψ SNI,PA,GO	[80,81]
Human Immunodeficiency Virus (HIV)-1 ^ψ	Acquired Immunodeficiency Syndrome (AIDS)	^ψ DC,BC,DD,SNI,GO,PA(KEGG)	[93]
<i>Helicobacter pylori</i> ^φ	Gastritis, peptic ulcer and gastric cancer	^φ CNA,DA	[83]
<i>Campylobacter jejuni</i> ^φ	Gastroenteritis	^φ DC,CCf,SNI,CNA,GO,EPI	[84]
<i>Treponema pallidum</i> ^φ	Syphilis	^φ SNI,CNA	[85]
<i>Mycoplasma pneumoniae</i> ^φ	Atypical pneumonia	^φ SNI	[86]
<i>Mycobacterium tuberculosis</i> ^{φ,ψ}	Tuberculosis	^φ DC,CC,DD,APL,CCf,ND,Str,SNI,CNA ^ψ DC,GO,PA(IntAct)	[87,96]

Table 1 continued

Pathogen	Disease	Mode of Network analysis	References
<i>Staphylococcus aureus</i> ^φ	Abscesses, Furuncles, Atopic dermatitis	^φ DC,BC,EPI	[44]
<i>Bacillus anthracis</i> ^ψ	Anthrax	^ψ DC,BC,CNA,GO	[94]
<i>Francisella tularensis</i> ^ψ	Pneumonia	^ψ DC,BC,CNA,GO	[94]
<i>Yersinia pestis</i> ^ψ	Pneumonic, septicemic, and bubonic plagues	^ψ DC,BC,APL,CNA,GO, PA	[94,95]
<i>Plasmodium falciparum</i> ^{φ,ψ}	Malaria	^φ DC,BC,CC,APL,CCf, Str,Ecc,Ra,Ass,WI, Cn,ND,GO,PA ^ψ SNI, GO	[99,100]

^φ - analyses of intra-pathogenic system, ^ψ - analyses of host-pathogen systems

DC-degree centrality, BC-betweenness centrality, DD-degree distribution, APL-average path length, CCf-clustering coefficient, CC-closeness centrality, ND-network diameter, NC-neighbor connectivity, Str-stress, Ecc-eccentricity, Ra-radiality, Ass-assortativity, WI-Weiner Index, Cn-centroid, DA-domain analysis; CNA-comparative network analysis, PA-pathway analysis, SNI-subnetwork identification, EPI-essential protein identification, GO-gene ontology analysis.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Figure 1: Schematic representation summarizing different computational approaches to analyze protein interaction networks
169x158mm (300 x 300 DPI)



Supplementary Table1: The number of intra-pathogenic and host-pathogen PPIs from different resources along with pathogens' genome size

Pathogen	Genome Size	Interactions		Resources
<i>Bacillus anthracis</i>	5.2 Mb	3036 ^φ	3020 ^ψ	PATRIC
<i>Campylobacter jejuni</i> NCTC 11168	1.6 Mb	24023 ^φ		PATRIC
<i>Clostridium botulinum</i>	3.9 Mb	10 ^φ	8 ^ψ	PATRIC
<i>Escherichia coli</i> O157:H7	5.6 Mb	3027 ^φ	12 ^ψ	PATRIC
<i>Francisella tularensis</i>	1.9 Mb	1312 ^φ	1312 ^ψ	PATRIC
<i>Helicobacter pylori</i>	1.7 Mb	2784 ^φ	4 ^ψ	PATRIC
<i>Listeria monocytogenes</i>	2.9 Mb	5 ^φ	5 ^ψ	PATRIC
<i>Mycobacterium tuberculosis</i>	4.4 Mb	8042 ^φ		[82]
<i>Mycoplasma pneumonia</i>	0.81 Mb	178 ^φ		[82]
<i>Pasteurella multocida</i>	2.4 Mb	12 ^ψ		HPIDB
<i>Salmonella enterica</i> Typhi	4.8Mb	10 ^φ	5 ^ψ	PATRIC
<i>Shigella flexneri</i>	4.6 Mb	191 ^φ	41 ^ψ	PATRIC
<i>Staphylococcus aureus</i>	2.9 Mb	16 ^φ	21 ^ψ	PATRIC
<i>Streptococcus pneumoniae</i> TIGR4	2.1 Mb	429 ^φ	12 ^ψ	PATRIC
<i>Treponema pallidum</i>	1.1 Mb	3649 ^φ		[82]
<i>Vibrio cholera</i> O1 biovar El Tor N16961	4 Mb	9 ^φ	1 ^ψ	PATRIC
<i>Yersinia pestis</i>	4.7 Mb	3948 ^φ	4018 ^ψ	PATRIC, HPIDB
<i>Epstein-Barr virus</i>	171.8 Kb	220 ^φ		BioGRID
<i>Hepatitis C virus</i>	9.6 Kb	111 ^φ		BioGRID
<i>Human immunodeficiency virus</i>	9.1 Kb	1195 ^φ		BioGRID
<i>Influenza A virus</i> (H1N1)	13.6 Kb		4067 ^ψ	HPIDB
<i>Kaposi's sarcoma-associated herpesvirus</i>	138 Kb	142 ^φ		BioGRID
<i>Severe acute respiratory syndrome-coronavirus</i>	29.8 Kb	65 ^φ		[82]
<i>Varicella zoster virus</i>	125Kb	173 ^φ		[82]
<i>Plasmodium falciparum</i>	22.9 Mb	4750 ^φ	367 ^ψ	[100]

φ and ψ represent the number of interactions for intra-pathogenic and host-pathogen systems, respectively

Supplementary Table1: The number of intra-pathogenic and host-pathogen PPIs from different resources along with pathogens' genome size

Pathogen	Genome Size	Interactions		Resources
<i>Bacillus anthracis</i>	5.2 Mb	3,036 ^φ	3020 ^ψ	PATRIC
<i>Campylobacter jejuni</i> NCTC 11168	1.6 Mb	24,023 ^φ		PATRIC
<i>Clostridium botulinum</i>	3.9 Mb	10 ^φ	8 ^ψ	PATRIC
<i>Escherichia coli</i> O157:H7	5.6 Mb	3,027 ^φ	12 ^ψ	PATRIC
<i>Francisella tularensis</i>	1.9 Mb	1,312 ^φ	1312 ^ψ	PATRIC
<i>Helicobacter pylori</i>	1.7 Mb	2,784 ^φ	4 ^ψ	PATRIC
<i>Listeria monocytogenes</i>	2.9 Mb	5 ^φ	5 ^ψ	PATRIC
<i>Mycobacterium tuberculosis</i>	4.4 Mb	8,042 ^φ		[82]
<i>Mycoplasma pneumonia</i>	0.81 Mb	178 ^φ		[82]
<i>Pasteurella multocida</i>	2.4 Mb	12 ^ψ		HPIDB
<i>Salmonella enterica</i> Typhi	4.8Mb	10 ^φ	5 ^ψ	PATRIC
<i>Shigella flexneri</i>	4.6 Mb	191 ^φ	41 ^ψ	PATRIC
<i>Staphylococcus aureus</i>	2.9 Mb	16 ^φ	21 ^ψ	PATRIC
<i>Streptococcus pneumoniae</i> TIGR4	2.1 Mb	429 ^φ	12 ^ψ	PATRIC
<i>Treponema pallidum</i>	1.1 Mb	3,649 ^φ		[82]
<i>Vibrio cholera</i> O1 biovar El Tor N16961	4 Mb	9 ^φ	1 ^ψ	PATRIC
<i>Yersinia pestis</i>	4.7 Mb	3,948 ^φ	4018 ^ψ	PATRIC, HPIDB
<i>Epstein–Barr virus</i>	171.8 Kb	220 ^φ		BioGRID
<i>Hepatitis C virus</i>	9.6 Kb	111 ^φ		BioGRID
<i>Human immunodeficiency virus</i>	9.1 Kb	1,195 ^φ		BioGRID
<i>Influenza A virus</i> (H1N1)	13.6 Kb		4067 ^ψ	HPIDB
<i>Kaposi's sarcoma-associated herpesvirus</i>	138 Kb	142 ^φ		BioGRID
<i>Severe acute respiratory syndrome-coronavirus</i>	29.8 Kb	65 ^φ		[82]
<i>Varicella zoster virus</i>	125Kb	173 ^φ		[82]
<i>Plasmodium falciparum</i>	22.9 Mb	4,750 ^φ	367 ^ψ	[100]

φ and ψ represent the number of interactions for intra-pathogenic and host-pathogen systems, respectively