

# SCIENTIFIC REPORTS



OPEN

## Understanding cancer complexome using networks, spectral graph theory and multilayer framework

Aparna Rai<sup>1</sup>, Priodyuti Pradhan<sup>2</sup>, Jyothi Nagraj<sup>3</sup>, K. Lohitesh<sup>3</sup>, Rajdeep Chowdhury<sup>3</sup> & Sarika Jalan<sup>1,2</sup>

Received: 31 March 2016  
Accepted: 15 December 2016  
Published: 03 February 2017

Cancer complexome comprises a heterogeneous and multifactorial milieu that varies in cytology, physiology, signaling mechanisms and response to therapy. The combined framework of network theory and spectral graph theory along with the multilayer analysis provides a comprehensive approach to analyze the proteomic data of seven different cancers, namely, breast, oral, ovarian, cervical, lung, colon and prostate. Our analysis demonstrates that the protein-protein interaction networks of the normal and the cancerous tissues associated with the seven cancers have overall similar structural and spectral properties. However, few of these properties implicate unsystematic changes from the normal to the disease networks depicting difference in the interactions and highlighting changes in the complexity of different cancers. Importantly, analysis of common proteins of all the cancer networks reveals few proteins namely the sensors, which not only occupy significant position in all the layers but also have direct involvement in causing cancer. The prediction and analysis of miRNAs targeting these sensor proteins hint towards the possible role of these proteins in tumorigenesis. This novel approach helps in understanding cancer at the fundamental level and provides a clue to develop promising and nascent concept of single drug therapy for multiple diseases as well as personalized medicine.

The post-genomic era aims to understand human health and diseases by investigating the role of proteomics and genomics, that involves macromolecules such as the proteins and nucleic acids (e.g. DNA, RNA, miRNA etc)<sup>1</sup>. Cancer being a multifactorial disease can be studied through these macro-molecules. To understand this complexome, there has been rapid advancements in both experimental and theoretical techniques in cancer diagnostic and screening<sup>2</sup>. These investigations indicate that all the cancers share a common pathogenic mechanism<sup>3</sup>. Much like Darwinian evolution, cancer cells acquire continuous heritable genetic variation by arrays of random mutation and go through the process of natural selection resulting in phenotypic diversity<sup>4,5</sup> like, differential gene expressions, alterations in cell regulation and control mechanisms, alteration in macromolecular interaction pathways, etc. These two fundamental processes in cancer cells provide them the capacity to have proliferative advantage and higher rate of survival than their neighboring cells<sup>4</sup> resulting in heterogeneous tumor formations<sup>6</sup>. This heterogeneity is found in both intra- and inter-tumor cell populations<sup>7</sup>. In addition, there are non-genetic factors that result in phenotypic diversity, e.g. epigenetic modifications, clinical diagnostic and therapeutic responses<sup>8,9</sup>. All these factors result to aberrations in various biological processes of the cancer cells and make cancer a complexome with no direct correspondence between the cancer and the normal tissues. These studies have remarkably improved our understanding of various factors associated with the cancer. However, even after billion dollars of investments<sup>10,11</sup> and extensive research, the major challenge lies in understanding the angiogenic and metastatic complexity<sup>12</sup>, modeling the disease at a global scale, drug target identification and co-evolving tumor cell<sup>13</sup>. These challenges forms the backbone of cancer systems biology. The research involving genetics at the molecular level have identified a number of susceptible genes responsible for the genesis of different types of cancers<sup>14</sup>. However, out of about 10% of the total cancer genes only 1% are known for their biological functions, indicating that the etiology of cancer is still not clear<sup>15</sup>. This demands for a holistic approach to understand cancer from a fundamental point of view. One such promising approach is to consider the cancer system as networks.

<sup>1</sup>Centre for Biosciences and Biomedical Engineering, Indian Institute of Technology Indore, Simrol, Indore, Madhya Pradesh 453552, India. <sup>2</sup>Complex Systems Lab, Discipline of Physics, Indian Institute of Technology Indore, Simrol, Indore, Madhya Pradesh 453552, India. <sup>3</sup>Department of Biological Sciences, Birla Institute of Technology and Science, Vidya Vihar, Pilani, Rajasthan 333031, India. Correspondence and requests for materials should be addressed to S.J. (email: sarikajalan9@gmail.com)

Studying cancer under network theory framework has already helped in understanding various constitutive properties of the system<sup>16–18</sup>. Various network studies pertaining to, epigenetic modifications, gene regulations, gene expressions, protein-protein interactions (PPI) have provided insights into the molecular mechanisms of the disease. Additionally, these network studies have helped in finding functionally important proteins as well as some of the missing pathways in cancer<sup>19–22</sup> providing a global understanding to biological processes and protein interactions<sup>23–27</sup>.

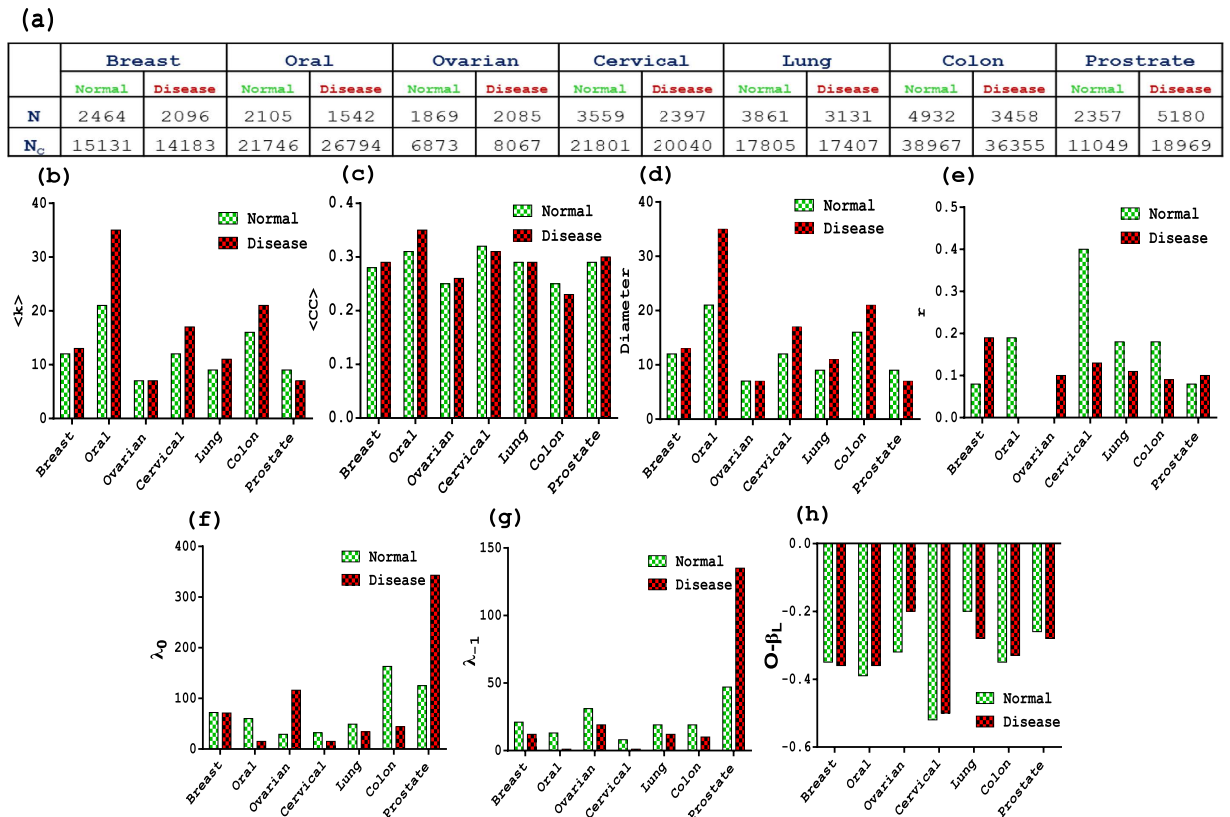
In this work, we consider seven most prevalent in human cancers namely, the breast, oral, ovarian, cervical, lung, colon and prostate. We analyze the protein-protein interactions among the normal and disease cells using the combined framework of network theory, spectral graph theory and the multilayer analysis to understand cancer development, progression and treatment response. The spectral graph theory has shown its remarkable success in uncovering the behavior of various complex systems<sup>27–34</sup> including biological systems corresponding to gene co-expression and PPI networks<sup>35–38</sup>. Further, implementing the multilayer analysis, we scale these seven cancers on the basis of the presence of common proteins into three different categories elaborately discussed in the methods section. The network approach of analyzing seven cancers not only demonstrates deviation in the complexity of all the normal and disease datasets from their corresponding random networks but also depicts changes in the complexity level between the normal and the disease states, contributing to understand cancer at the fundamental level. The multilayer framework highlights the proteins which are common in all the cancers and have structural importance in individual networks. Importantly, these common proteins also exhibit functional importance for occurrence of cancer revealed through pathway ontology and miRNA analysis. The framework paves a new way to the promising and nascent concept of single drug therapy<sup>39</sup> for multiple diseases as well as personalized medicine<sup>40</sup> in a time efficient and cost effective manner.

## Results and Discussions

**Properties of Complexome.** *Structural Properties.* We determine various structural properties of all the seven cancers for the normal and the disease states (Table S1 and S2). Additionally, we perform the comparative analysis of various properties of these networks with those of their corresponding random networks. We construct the corresponding Erdős Rényi (ER) random networks with the same  $N$  and average degree as of the PPI networks. The ER network only preserves the  $N$  and  $\langle k \rangle$ , but have interactions among pairs randomly chosen. Another random network model which we consider here is the configuration model. The configuration model additionally preserves the degree sequence of a network for the normal and the disease states. As indicated from Fig. 1a, all the disease datasets consist of different number of proteins as compared to their individual normal datasets. It is well implicated in various cancer proteome related studies that advantageous genetics and natural selection<sup>6</sup> lead to mutations in the proteins and impaired functions, causing different count of proteins participating in cancer networks<sup>41</sup>.

To get an overview of the structural organization of the protein interactions, we first examine the global structural properties of the normal and the disease networks. First such property is the number of connections. As indicated in Fig. 1b, the overall average degree is higher in the disease networks except for the ovarian and prostate cancer, indicating more tendency of proteins to interact with each other in these cancers. Other PPI studies have also reported that most of the cancer proteins exhibit a higher binding tendency to interact with other proteins due to favorable mutations at binding sites<sup>21</sup>. Hence, the observation of higher number of interactions in the disease state is not surprising. In the seven cancers considered here, the disease network of oral cancer exhibits relatively a higher number of connections than that of all the other networks suggesting that for this particular cancer, proteins have more affinity to interact among themselves. The observation is already reported that they have large molecular organizations facilitating maximum interactions with other proteins<sup>22</sup>. Further, the degree distribution for all these PPI networks follow power law behavior which is also observed for many other biological systems<sup>16</sup>. More intriguingly, the degree distribution of the all the networks follow two distinct fitting scales, i.e. two power law for all the networks (Fig. S1). Many real world networks have been reported to devoid of a perfect power law for the entire range of the degree<sup>42,43</sup>. The existence of two power law implicates that the network is robust not only by the inclusion of new proteins and interactions by the hub proteins but also by the contribution of new or altered interactions of existing proteins<sup>42</sup>.

Second such global structural property is diameter ( $D$ ) of the network which captures the spread of information in a system<sup>44</sup>. Diameter of a network is the longest of the shortest paths between all the pair of nodes of the network<sup>44</sup>. Different networks having same values of  $N$  and  $N_C$  (number of connections) can have different diameter<sup>45</sup>. Similarly, networks having different values of  $N$  and  $N_C$  can yield the same diameter as diameter is decided by the manner in which connections are distributed in a network. A small diameter is known for the faster signaling in the system<sup>45</sup>. Based on the diameter analysis, we have two results; first the diameter of PPI networks is larger than that of their corresponding ER and configuration random models (Table S2). Second, the disease PPI networks display relatively smaller diameter as compared to the corresponding normal ones. The smaller diameter in the disease state may be due to both the smaller size as well as large average degree of the disease networks (Table S2). Since, diameter has been emphasized to shed light on information propagation in a network<sup>45</sup>, a small diameter of the disease network indicates a faster propagation of signals in terms of less number of the intermediate paths involved. It is already known that disorders in cancer associated proteins promote the adaptability of faster communication in many major cancer related cellular signaling processes by up regulation or down regulation of pathways leading to uncontrolled cell proliferation<sup>46,47</sup>. Short path for PPI networks can also be considered with respect to time as few molecules such as miRNAs facilitating the disease proteins to regulate their expressions by mediating inter-cellular cell signaling in cancer cells which further lead to faster information transduction within the cell<sup>48</sup>. Note that, here networks are constructed by taking PPI from the same database (STRING) for both the normal and disease states and hence these networks do not capture dynamical behavior of the PPIs. A possible way of capturing the dynamical nature of PPIs is to consider gene-co-expression



**Figure 1. Different structural and spectral properties.** The table (a) summaries number of nodes ( $N$ ) and number of connections ( $N_c$ ) of all PPI networks. The graph represents (b) the average degree ( $k$ ), (c) average clustering coefficient ( $C$ ), (d) diameter  $D$ , (e) degree-degree coefficient  $r$ , degenerative eigenvalues: (f)  $\lambda_{-1}$  and (g)  $\lambda_0$ , and (g) betweenness-overlap correlation ( $O-\beta_L$ ) for the real and their corresponding random networks for normal and disease datasets. All the cancers show similar statistics for  $\langle CC \rangle$ , diameter,  $\lambda_{-1}$ ,  $\lambda_0$ , except in prostate cancer. There is no significant comparison in the values of  $\langle CC \rangle$  and  $r$  between normal and corresponding disease networks.

networks<sup>49,50</sup> etc., but gene-co-expression networks also have their own limitations and drawbacks such as small sample size<sup>51</sup>.

Third global property indicating the interaction complexity of networks is the degree-degree correlation coefficient ( $r$ ) which indicates tendency of nodes to connect with (dis)similar degree nodes<sup>52</sup>. Almost all the networks possess positive  $r$  values except for the oral disease and the ovarian normal dataset. Since the  $r$  values of these two datasets is very close to zero, we can consider them to be neutral. Overall, both the disease and the normal PPI networks display positive or near to zero values of  $r$ , whereas the corresponding ER networks have  $r$  value zero. This is not surprising as the networks with the same average degree and size may still differ significantly in various network features based on their nature of interactions. In the ER random network, the nodes are randomly connected and  $r$  value of the network is determined by degrees of the interacting nodes. Further, in order to see whether changes in the degree-degree correlations of different networks are arising due to the change in the degree sequence, we compare all the PPI networks with their corresponding configuration networks. The  $r$  value of the fourteen configuration networks turns out to be overall dissortative in nature. Therefore, some normal networks are more assortative than the disease networks while, vice-versa behavior is seen for some cancers (Fig. 1e). Since, ER random and configuration model depict similar behavior of all the fourteen networks, the changes in the  $r$  values of the networks constructed from the datasets have the following implications. There are changes in the interaction patterns of the disease from the normal networks as the change in  $r$  is not arising due to changes in the degree sequence or data size. This further implicates that there exists varying complexity in the underlying system. Complexity of molecular associations in different normal tissues can be understood as there are different proteins and pathways that define cell proteome and various sub-tissue types which are functioning for particular tissue fate. Similarly, Cancer genetics and histopathology comprises of -intra and -inter tumor heterogeneity and its metastasis may harbor overall divergent biological behavior. In other words, different cancers may have altogether different cancer genetics and histopathology in their tumors<sup>7</sup>. Altogether, the unsystematic changes in the values of  $r$  from the normal to the disease networks indicate the inherent complexity of the two states.

Analysis of above global properties of the network indicate deviations of the real networks from their random counterparts as well as reflects overall change in the complexity of interaction patterns of the normal and the disease networks. We further investigate properties providing understanding of the local interactions in the network. One such property is the clustering coefficient of nodes in a network which suggests how the neighbors of

a node are connected<sup>45,53</sup>. The average clustering coefficient ( $\langle CC \rangle$ ) depicts an unsystematic change in its values in the disease and their normal networks, i.e. some disease datasets have more  $\langle CC \rangle$ , while for some datasets normal networks have more  $\langle CC \rangle$  (Fig. 1c). Different values of  $\langle CC \rangle$  further indicates changes in the interaction patterns in the normal and disease states. But what is noteworthy is that these networks have much higher  $\langle CC \rangle$  than those of their corresponding ER and configuration networks, indicating the abundance of cliques of order three<sup>54</sup>. Cliques indicate the preserved interactions in the networks and are believed to be conserved during the process of evolution<sup>55</sup>. Further, these structures are also considered to be building blocks of a network, making the underlying system more robust<sup>56</sup> and stable<sup>57</sup> and cancer as a system is reported to be robust against both the targeted chemotherapy and the hazardous environment<sup>58,59</sup>. Another property revealing structural complexity of a network is the correlation between degree and clustering coefficient ( $k - CC$ ) for the nodes in a network. The  $k - CC$  correlations of all the disease and normal networks manifest overall negative value (Fig. S2), as also exhibited by many other biological systems reflecting the presence of hierarchical structure in the system<sup>16</sup>. The presence of hierarchical structures is an indication of highly clustered neighborhoods consisting of sparsely connected nodes communicating through hubs and functional modules in the network<sup>54</sup>.

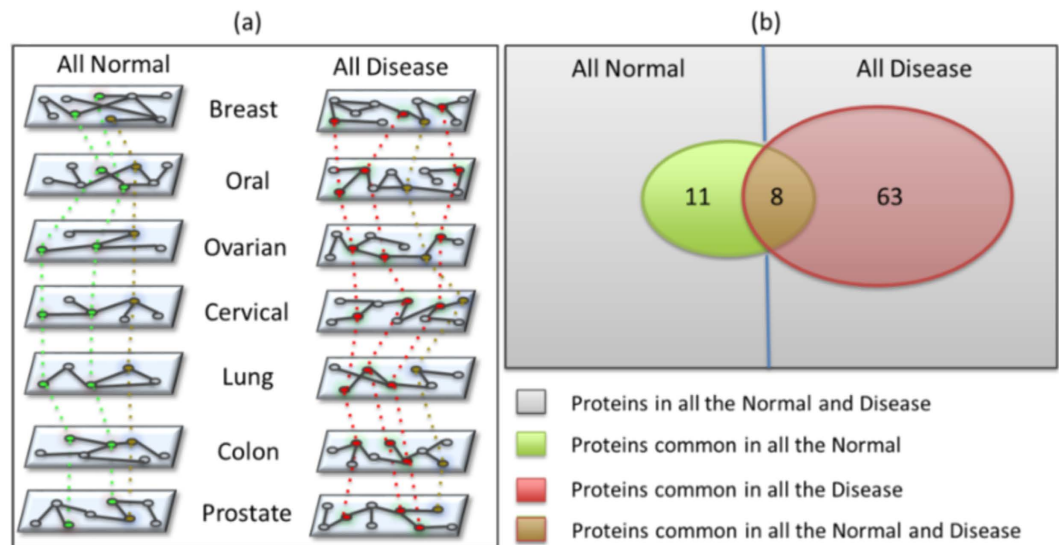
All the structural properties discussed above reveal overall similar behavior for almost all the normal and the disease states as well as indicate existence of complexity in interaction pattern of both the networks. Further, the variation in the values of  $r$  and  $\langle CC \rangle$  from the normal to the disease, i.e. in some of the datasets normal having more values of  $r$  and  $\langle CC \rangle$  than those of disease and vice versa, which reflect an unsystematic change in the interactions from the normal to the disease networks. To have an in-depth analysis of proteins or factors causing changes in the interaction patterns in the PPI networks considered here, we explore the spectral properties of these networks.

**Spectral properties.** The eigenvalue distributions of all the normal and disease networks depict triangular shape (Fig. S3) as observed for many other biological networks<sup>60</sup>. Furthermore, the spectra exhibit a very high degeneracy at the zero eigenvalue for all the networks as compared to that of their corresponding ER random networks. The corresponding configuration models exhibit a high degeneracy at the zero eigenvalue which indicates that not only a particular degree sequence but also the nature by which these protein-protein interactions contribute on occurrence of the high degeneracy at the zero eigenvalues in the real networks. Since number of zero eigenvalues in the adjacency matrix is directly related with the complete and partial node duplicates in the underlying network<sup>61</sup>, a very high value of  $\lambda_0$  degeneracy indicates occurrence of node duplication in these PPI networks. The duplicate nodes are the ones which shares the same neighbors in a given network. Here we consider the nodes which are complete duplicates, that is, these nodes have exactly the same neighbors. There are partial duplicate nodes also in the network which do not have exactly the same neighbors but possess few uncommon neighbors too. Finding partial duplicate nodes in a network is computationally very exhaustive and hence here we only concentrate on complete duplicate nodes. These duplicate nodes are known to be important during the evolution<sup>62,63</sup> and hence occurrence of node duplication in the normal networks is not surprising as it indicates the evolutionary processes over the years. Interestingly, duplication is also known to be one of the important factor in promoting cancer and contribute in evolving the normal cell to the disease state<sup>64,65</sup>. But, what is interesting is that despite a very high  $\lambda_0$  degeneracy in the disease networks, indicating a very high number of exact and partial duplicates, most of the duplicate nodes in the disease PPI networks are different from those of the corresponding normal PPI networks (Fig. S4). This observation suggests that the genetic mutations leading to abrupt transformation of a normal cell into a cancerous cell may have caused incurrence of new proteins to perform similar functions and thus resulting in the duplicate nodes in the disease state. Moreover, there is a different number of the duplicate nodes in different cancers which can be understood in terms of independent adaptation of each cancer genome arising due to independent heterogeneity and natural selection<sup>66</sup>. Further, both the disease and the normal networks display the degeneracy at minus one eigenvalues ( $\lambda_{-1}$ ) whereas  $\lambda_{-1}$  is absent in both the corresponding ER and the configuration networks. It indicates that the PPI networks may contain complete sub-graphs. Complete sub-graphs are known to be the building blocks of a network further making the network robust<sup>56</sup> and stable<sup>57</sup>. It is known that architecture of PPI network is made up of sub-networks of metabolic cycles and pathways which play important role in constituting PPI networks<sup>67,68</sup>.

One striking observation is that both the disease and the normal datasets of prostate cancer result in different network properties than those of the other cancer datasets. For instance, some of the properties analyzed here, like the average degree ( $\langle k \rangle$ ), diameter and minus one degenerate eigenvalues of the prostate cancer are exception to that of the other networks. The disease network of prostate has much higher value of zero eigenvalues than the normal, whereas for other pairs the vice versa behavior is observed with an exception of ovarian cancer. For other properties such as the average clustering coefficient and degree-degree correlation coefficient, there is unsystematic changes from normal to disease. Different behavior of prostate cancer may be arising due to lack of availability of complete knowledgebase of proteomic interactions for prostate cancer. Further, it has been reported that prostate cancer is very lately diagnosed<sup>69</sup> and thus, the altered network properties of prostate cancer may also suggest the significance of independent cancer development processes in this cancer.

**Multilayer analysis.** All the structural and spectral properties reveal similar behavior for almost all the seven normal and the disease networks such as high value of  $\langle CC \rangle$ , smaller diameter as compared to their corresponding random networks, non-negative  $r$  values, negative  $k - CC$  correlations and the triangular distribution of the eigenvalues with very few exceptions for each which have been discussed in the above section. However, the disease state differ from their normal counterpart in terms of  $r$  values and  $\langle CC \rangle$ , suggesting difference in the interaction patterns between the disease and the normal networks. To get insight into the proteins responsible for the changes in the interaction patterns from the normal to the disease, which may also be crucial in making a normal tissue to the cancerous one, we enlist the common proteins in all the normal as well as disease datasets (Fig. 2)





**Figure 2. Multilayer analysis.** (a) Schematic diagram showing the construction of multilayer network where each normal and disease network of the seven cancers are represented as layers leading to normal and disease multilayer networks respectively. The dotted lines represent the common proteins considered from each of dataset, the red, green and blue circles represent common proteins in all (i) the disease, (ii) the normal datasets and both the normal and disease datasets (union of (i) and (ii)), respectively. After extraction of these common proteins, their interaction partners are taken from individual datasets. (b) The Venn diagram of common proteins depicting the number of proteins common in all the normal and disease dataset, respectively.

as explained in the multilayer framework in the methods section. There are 63 proteins which are common in all the disease networks. If a protein appears in all the disease dataset, it is enlisted in the common proteins list. We investigate the pathways involved by considering the interactions of these common proteins as well as their structural importance in the networks. There are 19 proteins which appear common in all the seven normal datasets and 71 in all the seven disease ones. Out of these, 8 proteins appear common in all the normal and all the disease datasets (Fig. 2b). The common proteins occur in different cancers as the tumor cells share similar cellular environment biologically i.e. uncontrolled cell division, cell proliferation, metastasis etc. There are some reports on the occurrence of common protein markers in different cancers<sup>70</sup> which suggest that there might be similar features associated with these common proteins present in different cancers. We investigate the network properties of these common proteins for all the disease states to find their importance in the network architecture and also study their biological functions.

First, we extract the interacting partners of all these common proteins from the individual disease networks (Table S5 and S6). We find that though the proteins are common in these PPI networks, some of the interacting partners of these common proteins are different in individual networks suggesting addition or deletion of proteins due to mutations caused in each cancer. Thereafter, we enlist the number of interactions among the 63 common proteins (referred as IN connections) and the interactions outside the 63 proteins (referred as OUT connections) (Table S6). These proteins have much higher number of OUT connections than the average degree of the corresponding network (about two fold of average degree), suggesting their significant contribution in the overall network connectivity. We further analyze interaction properties of these subtractive PPI networks. We determine the  $\langle CC \rangle$  of these 63 nodes and compare it with that of the whole network for all the seven disease datasets. The  $\langle CC \rangle$  of 63 nodes in each disease is much higher (nearly twice) as compared to the corresponding whole networks. A high  $CC$  of 63 proteins indicates accountability of these proteins for higher  $CC$  of the whole network as well as existence of modular structures in the network. To have a broader understanding into the organization of these 63 proteins, we perform molecular and pathway ontology of these common proteins. To do this, we retrieve protein sequences of 63 proteins from UniProtKB and direct them to Reactome pathway browser<sup>71</sup>. We find that, among the 63 common proteins, many proteins collectively play significant role in few of the important pathways such as Vascular endothelial growth factor (VEGF) signaling, PIP3 activation, VEGFA-VEGFR2 Pathway, Signaling by PDGF, Signaling by NGF etc (Table 1). All of these pathways are reported to have constitutive role in different cancers (detailed description in Supplementary Material).

Furthermore, we investigate other structural properties such as  $k - \beta_c$  (degree-betweenness centrality) correlation, weak ties analysis etc of these 63 proteins. The analysis reveals proteins contributing to the occurrence of the disease.

**$k - \beta_c$  correlation and weak ties analysis.**  *$k - \beta_c$  correlation.* We analyze the correlation between the degree and betweenness centrality for all the disease networks and highlight the 63 proteins as shown in Fig. 3. All the disease networks exhibit overall positive  $k - \beta_c$  correlation as also observed for many other complex systems<sup>16</sup>. Although few networks depicts data being clustered, i.e., data-points are localized in several regions, the

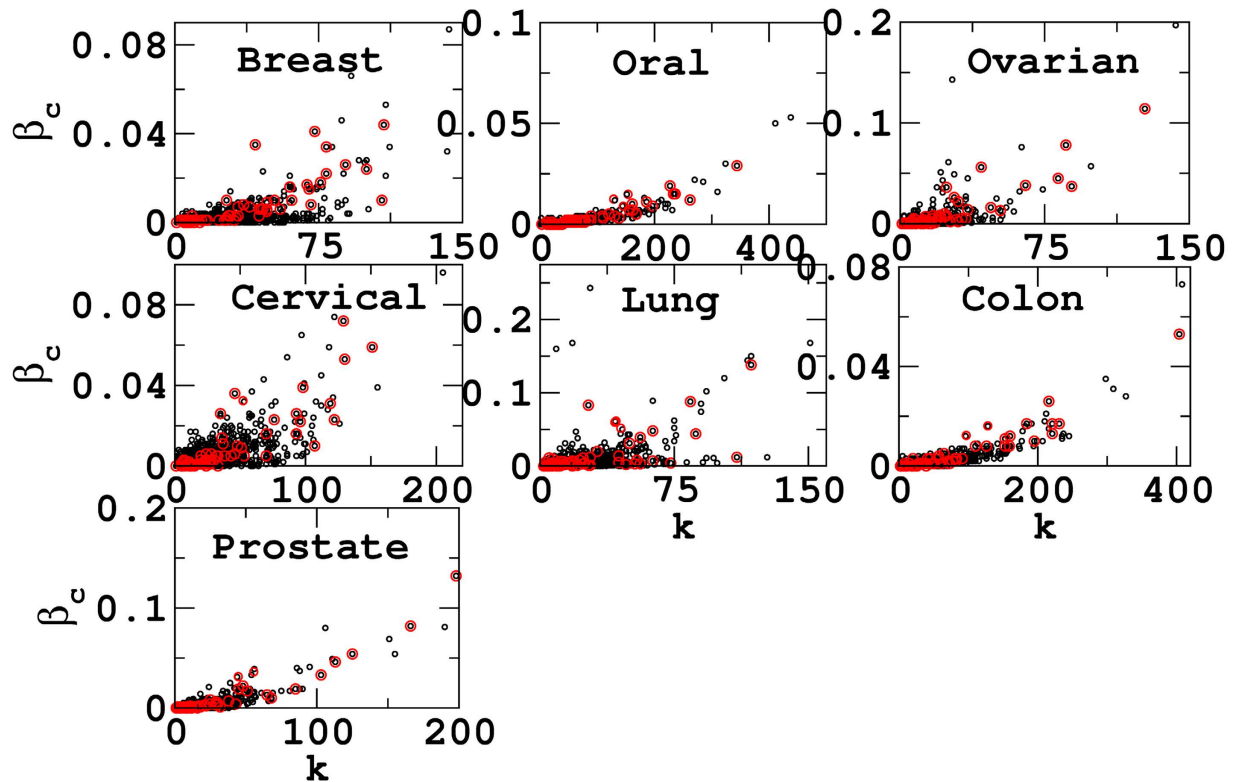
Sr.	Pathway	Proteins involved
1.	Signaling by Vascular endothelial growth factor (VEGF)	IQGAP1, FN1, PTK2, AKT1, FGFR2, MAPK1, VEGFA, CTNNA1, CTNNB1, CAV1
2.	Signaling by Stem cell factor receptor (SCF-KIT)	IQGAP1, FN1, STAT1, GSK3B, PTK2, AKT1, FGFR2, MAPK1, PTEN
3.	VEGFA-VEGFR2 (VEGF family receptors) Pathway	IQGAP1, FN1, PTK2, AKT1, FGFR2, MAPK1, VEGFA, CTNNA1, CTNNB1, CAV1
4.	Signaling by epidermal growth factor receptor 4 (ERBB4)	IQGAP1, FN1, ESR1, GSK3B, PTK2, AKT1, FGFR2, MAPK1, PTEN
5.	Protein kinase B (AKT) signaling	GSK3B, AKT1, FGFR2, PTEN
6.	Cellular responses to stress	HSPA4, PRDX5, EP300, GSK3B, SOD2, MAPK1, NBN, VEGFA, HSPA5, PRDX2, PRDX1
7.	Signaling by Platelet-derived growth factor (PDGF)	IQGAP1, FN1, STAT1, GSK3B, PTK2, AKT1, FGFR2, MAPK1, PTEN
8.	Downstream signaling of activated FGFR2	IQGAP1, FN1, GSK3B, PTK2, AKT1, MAPK1, FGFR2, PTEN
9.	Signaling by Nerve Growth Factor (NGF)	IQGAP1, FN1, GSK3B, PTK2, AKT1, ARHGDI, MAPK1, FGFR2, PTEN, RTN4
10.	Signaling by Rho GTPases	IQGAP1, BIRC5, CDH1, PTK2, ARHGDI, MAPK1, SFN, CTNNA1, CTNNB1, CTTN
11.	Axon guidance	IQGAP1, FN1, GSK3B, PTK2, MMP2, FGFR2, MAPK1, VEGFA, CFL1
12.	Innate Immune System	IQGAP1, FN1, EP300, GSK3B, PTK2, AKT1, FGFR2, MAPK1, PYCARD, PTEN, CFL1, CTNNB1
13.	Signal Transduction	IQGAP1, FN1, STAT1, EP300, GSK3B, ARHGDI, SFN, CTNNA1, CTNNB1, CAV1, BIRC5, ESR1, CDH1, PTK2, AKT1, MAPK1, FGFR2, VEGFA, PTEN, RTN4, CTTN
14.	Metabolism of proteins	PDIA3, LMNA, PABPC1, BRCA1, MMP2, HSPA5, MUC1, CTNNB1, RPS3, PML

**Table 1. Molecular and pathway ontology of 63 common proteins.** Set of proteins are involved in particular cellular pathway having major role in occurrence of different types of cancers.

overall correlation behavior of these networks are positive. Further, there are few nodes which despite of having less number of interactions (low  $k$ ), participate in a large number of pathways calculated through the betweenness centrality (high  $\beta_c$ ). In biological context, these proteins may be important as betweenness measures the ways in which signals can pass through the interaction network. If a protein having a high  $\beta_c$  has low value of  $k$ , it depicts the participation of that protein in many pathways and connecting different functional modules. We find that there are many proteins having high  $\beta_c$  and low  $k$  in the individual disease datasets. However, here we consider only 63 disease common proteins which fall in this regime. Among 63 disease common proteins, six proteins possess remarkably high  $\beta_c$  and low  $k$  values in four (breast, ovarian, cervical and lung) of the seven disease networks. These proteins are namely MUC1, STAT1, SOD2, MAPK1, HSPA4 and HSPA5. In the other three disease networks (oral, colon and prostate) these proteins possess high  $\beta_c$ , but comparative to other four networks they do not have very low  $k$ . It is crucial to note here that these six proteins are not the hub proteins. The hub proteins have already been reported in carrying out many necessary and housekeeping functions in the cell<sup>72</sup>, and the list of significantly very high degree proteins or hub proteins can be found in Table S3. Let us focus on the six proteins revealed through  $\beta_c$  and  $k$  analysis having structural importance in the network architecture. We look for the biological functions of these proteins and find that these proteins are well implicated in cancers. All these six proteins, are involved in the anti-apoptotic pathways<sup>73</sup>. Particularly, five proteins, namely MUC1, STAT1, SOD2, MAPK1 and HSPA5 are additionally responsible to aggravate metastasis<sup>74</sup> and play a key role in tumor progression by acting as angiogenic regulators<sup>75</sup>. Further, three of these proteins STAT1, MAPK1 and HSPA5 also induce cell proliferation<sup>76</sup>. Thus, the proteins having high  $\beta_c$  and low degree exhibit significant involvement in cancer related activities. The detailed functional properties are summarized in the Supplementary Material.

After investigating the nodes which are important in various pathways, we direct our attention towards finding important links or edges in the network connecting the 63 proteins, through Granovetter's 'weak ties hypothesis'.

**Weak ties analysis.** As defined in the methods section, weak ties are the links with their end nodes having very less number of common neighbors<sup>77</sup>. This analysis, motivated from the social network research, highlights importance of an edge in a network through its edge-betweenness centrality<sup>78</sup>. An inverse relation between the strength of a tie and overlap of the neighbors of nodes at both the ends indicate an existence of weak ties in the network. If a link has a high betweenness centrality ( $\beta_L$ ), it is known to be stronger as it helps in connecting different modules in the network. These weak ties are cited to be important in connecting different communities<sup>79</sup>. For PPI networks considered here, we find negative  $O-\beta_L$  correlation for all the normal and disease networks which suggests the presence of weak ties in these networks. This is in line of the earlier observation that the PPI networks are known to compose of different metabolic cycles and biological pathways. A protein involved in particular pathway plays role in regulating other pathways as well, termed as cross talk between pathways<sup>67,68</sup>.



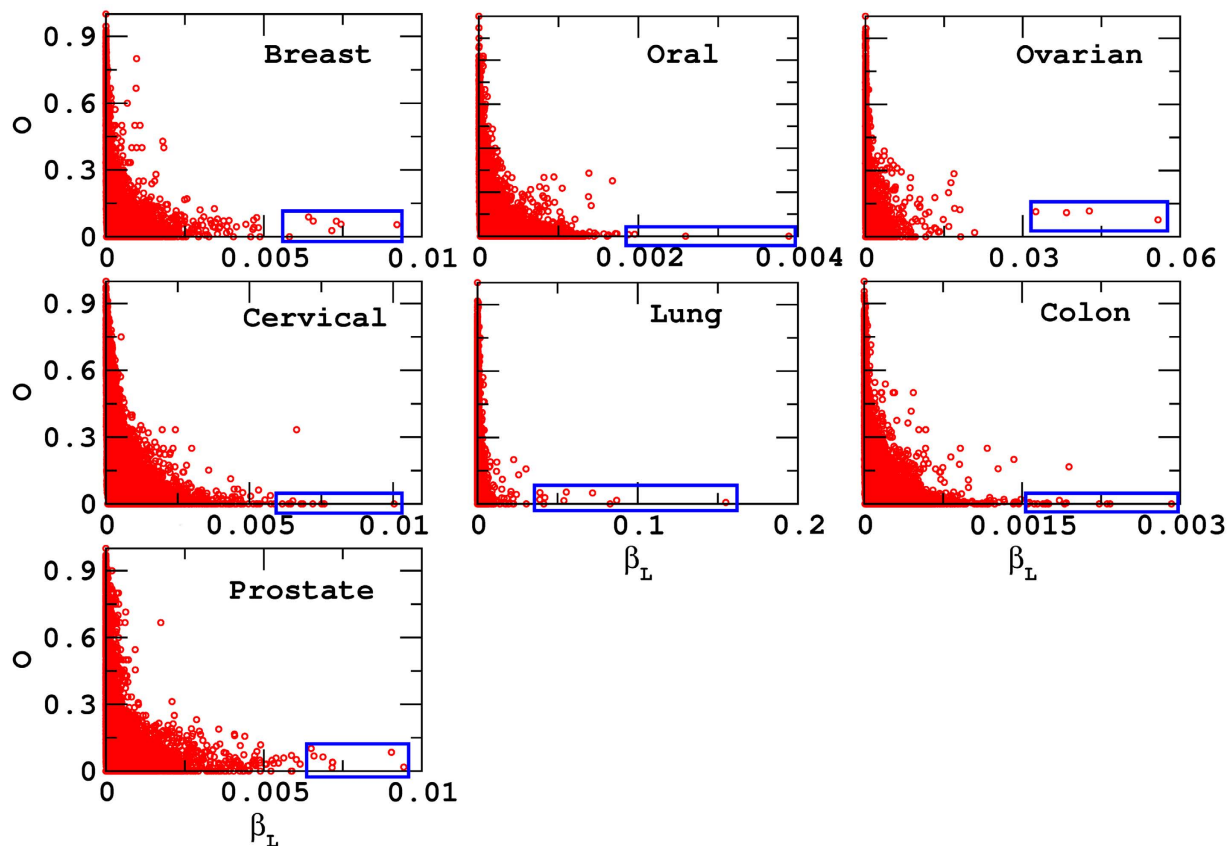
**Figure 3.  $k - \beta_c$  correlation.** The  $k - \beta_c$  correlation for all the disease networks reveal positive correlation (black circles). The red circles depict the  $k - \beta_c$  correlation for 63 disease common networks.

The weak ties analysis reveals 122 proteins (61 pairs) for all the disease networks together (Fig. 4 blue box). Among the 63 disease common proteins, ten proteins possess the properties of weak ties. Of these ten proteins, five proteins namely MUC1, SOD2, MAPK1, HSPA4 and HSPA5 are among the six proteins which we have already listed for their importance in possessing the property of high  $\beta_c$  and low  $k$ . This is not surprising as the nodes having high betweenness and low  $k$  are highly probable of having less overlap and thus the link betweenness of those nodes become high. The proteins having both the weak ties as well as high  $\beta_c$  low  $k$  properties (Fig. 5) indicates participation of these proteins in many pathways in the network as well as their significant role in causing cancer which we discuss in the following section.

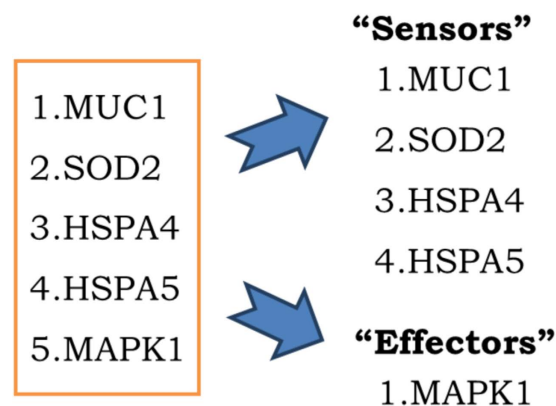
We discern the functional characterization of the proteins revealed in the above analysis based on their sub-cellular locations i.e. sensors and effectors. The sensors and effectors are widely characterized in direct or indirect involvement of a protein in cancer biology. Further, it is reported that post transcription regulators such as non-coding RNAs, particularly the miRNAs effectively regulate the expression of sensors<sup>80</sup>. miRNAs are a class of short non-coding RNAs with post transcription regulatory functions. Here, we study the role of miRNAs to understand the regulation of these proteins at the transcription level.

**Functional role and miRNA analysis of important proteins.** We find that out of five proteins participating in the weak ties as well as having high  $\beta_c$  low  $k$  properties, four proteins MUC1, SOD2, HSPA4 and HSPA5 are under the category of sensors and one protein MAPK1 is categorized under effectors (Fig. 5). The proteins marked under ‘sensors’ category are primarily upstream components of intra-cellular signaling cascades, altered expression of which may lead to downstream activities in the tumor milieu<sup>81</sup>. The protein under ‘effector’ category is downstream molecule which is often implicated but is not exclusive to cancer and therefore, for elaborative studies, we only discuss sensors (Fig. 5).

**Functional role.** We first scan the significant interactions of these proteins from the STRING database, which is based on probabilistic confidence score ( $>0.50$ ). The associations in STRING are based on high throughput experimental data, thorough search of the databases and predictions based on genomic context analysis. Thereafter, utilizing and incorporating information of interacting partners by KEGG pathway analysis, we find that all these sensor proteins HSPA4, HSPA5, MUC1 and SOD2 have very important interacting proteins e.g. ESR1, ErBb2, UBC, OS9 etc that have significant and specific involvement in the pathways participating in proliferation and migration of cancerous cells (Fig. 6(A)). Moreover, even after removing lower probabilistic interactions, these proteins possess some common neighbors viz. TP53, HSP90AA1, ESR1 and UBC which illustrates that these proteins are interlinked to each other and suppressing the expression of any of these may lead to the blockage of the cancer associated pathways (Fig. 6(A) red boxes). Further, we look into the miRNA molecules associated with these proteins that can help regulate the expression of these sensor proteins.



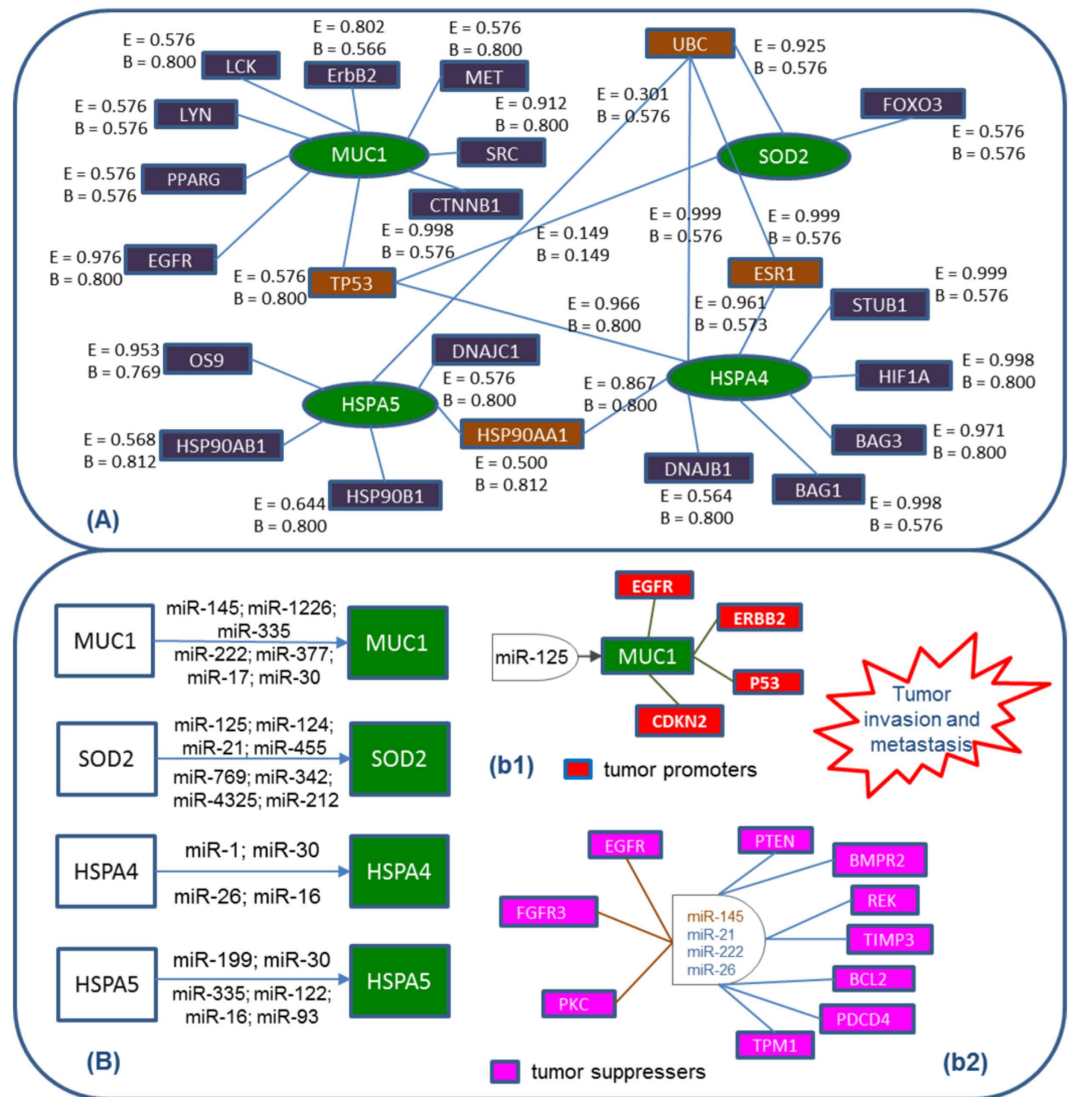
**Figure 4. Weak ties analysis.** The  $O-\beta_L$  analysis for all the disease networks reveal negative correlation (red circles). We highlight the edges (blue box) having high  $\beta_L$  and low  $O$  in all the disease and find the presence of 63 common proteins in those edges.



**Figure 5. Proteins with high  $\beta_c$ -lowk and weak ties.** MUC1, SOD2, HSPA4 and HSPA5 are sensors having role as upstream components of intra-cellular signaling cascades whereas MAPK1 is a downstream molecules classified here as effector.

**miRNA analysis.** Recent studies have shown that the expression of miRNAs is de-regulated in cancer progression, tumor invasion, metastasis, and subsequent chemoresistance<sup>82,83</sup>. For the four sensor proteins, we extract the list of experimentally validated miRNAs from Sanger Institutes miRTarBase database regulating the sensors proteins<sup>84</sup>. These miRNAs are then filtered for their role in regulating both the sensor proteins as well as their interacting partners (Fig. 6(B)). Here, we present example of few miRNAs corresponding to each sensor protein to demonstrate how miRNAs can be used to find out essential protein biomarkers and their downstream pathway roles. For instance, based on our studies, we find miR-125b as a probable miRNA regulating both MUC-1 and its interacting partners like, EGFR, ERBB2, CDKN2A which are potent tumor promoters (Fig. 6(B-b1)). Interestingly, KEGG pathway analysis, depicting various miRNAs de-regulated in cancer, reports that miR125 is





**Figure 6. Sensor proteins and their miRNAs associations.** (A) Depicts the interactions (blue boxes) of the sensor proteins (green ellipse) from STRING database. Red boxes are the common neighbors of the four sensor proteins. (B) Depicts the miRNAs associated with the sensor proteins. miRNA in (b1) is a positive regulator of MUC-1 leading to down regulation of tumor promoters (red boxes) while, other miRNAs in (b2) participate in up-regulation of MUC1, SOD2 and HSPA4 and down-regulate the activities of tumor suppressors (pink boxes).

down-regulated in various cancers which provides strong indication to the consequences of over-expression of MUC-1 in cancer (Fig. S6). Similarly, miR145 is predicted to have potential binding affinity for MUC1 which is suppressed in many cancers and is reported to have targets like, EGFR, FGFR3 and PKC (Fig. 6(B-b2)) which have well established role in tumorigenesis<sup>85</sup>. Further, miRNAs like, miR-21, miR-26 and miR-222 are up-regulated in various cancers and found to regulate MUC1, SOD2 and HSPA4. Existing literature suggests that these miRNAs also down-regulate tumor suppressors like, PTEN, BMPRII, RECK, TIMP3, BCL2, PDCD4, TPMI thereby promoting cancer (Fig. 6(B-b2)). Further, to have a complete idea about miRNA-mediated regulation, we calculate the probabilistic distribution of proteins regulated by a given miRNA which also controls the expression of sensor proteins and study the role these miRNAs play in regulating other proteins. It is revealed that the MAPKinase family is highly probable target (Fig. S9) since they are simultaneously being regulated by majority of miRNA which regulate sensor proteins. The implication of this investigation is that the proteins of the particular signaling pathway is highly important in cancer dataset under study and it can be chosen as a suitable target to be looked upon after miRNA inhibition. Apart from the MAPKs, other common targets are HSPs, B-catenin, PI3K/Akt, Mucin family, which is based on the probability scores alone (Table S10). In all, the data indicates the merits of using network theory to predict plausible proteins regulating a range of downstream targets. However, experimental validation is essential for a concrete conclusion.

## Conclusion

We analyze the protein-protein interaction networks of the normal and the disease conditions of seven different cancers under the combined framework of the spectral graph theory, network theory along with the multilayer analysis. The analysis exhibits overall similar behavior of various structural properties among the normal and the disease states (with few exceptions) such as a high clustering coefficient, small diameter, negative  $k - CC$  correlation and positive degree-degree correlation. Further, these properties exhibit significant deviations from those of their corresponding random networks. While, high  $\langle CC \rangle$  and negative  $k - CC$  correlation depicts complexity in the underlying system, change in the values of  $r$  and  $\langle CC \rangle$  from the normal to the disease states signify changes in the interaction patterns between the datasets. Further, the  $r$  and  $\langle CC \rangle$  values exhibit unsystematic changes from the normal to the disease networks, i.e. while some of the cancer networks have higher values of  $r$  and  $\langle CC \rangle$  than their normal counterpart, some have lower values as compared to the normal depicting inherent complexity in the normal and the disease networks.

Furthermore, to have a deeper insight into the complexity of the normal and the disease system and have in-depth analysis of the factors leading to changes in the interaction pattern in the networks, we analyze spectra of the cancer and normal networks, as well as compare them with those of their respective random models. We find that there is a high degeneracy at the zero and the minus one eigenvalue. The zero degeneracy is directly related with the number of duplicate nodes in the network and occurrence of different duplicate nodes in the normal and disease states suggest evolutionary changes from the normal to the disease state. Whereas, a high degeneracy at minus one eigenvalue suggests abundance of complete sub-graphs in the networks which may be important for proper functioning of the underlying system. Moreover, difference in the behavior of prostate cancer for various structural and spectral properties than those of the other cancer datasets may be due to incomplete knowledge-base of proteomic interactions or may have due to independent cancer development processes for this cancer.

Next, the multilayer framework reveals 63 common proteins among all the disease datasets. These 63 proteins show much higher  $\langle CC \rangle$  as compared to that of the whole networks. Further, the functional analysis of these 63 proteins through pathway ontology reflect their involvement in important cancer related pathways such as Vascular endothelial growth factor (VEGF) signaling, PIP3 activation, VEGFA-VEGFR2 pathway, signaling by PDGF, signaling by NGF etc.

Other network properties investigated for the 63 proteins common to all the diseases namely the  $k - \beta_i$  correlation and weak ties ( $O - \beta_i$ ) analysis highlights few proteins and their links that are structurally important in the network. These proteins are also found to have functional significance responsible for the occurrence of cancer. These proteins are further categorized into sensors and effectors. The sensors having primary role in contributing changes in the tumor<sup>81</sup>, are studied for post transcription regulators such as the miRNAs as they effectively regulate the expression of sensor proteins. Out of 5 common proteins which possess structural importance in the individual networks, 4 are sensor proteins. The miRNA study of these sensor proteins reveals their involvement in tumor invasion and metastasis thereby suggesting their role in progression of the cancer. Further experimental validations of these miRNAs can help in making corresponding proteins as potential drug targets.

All these results based on rigorous analysis using sophisticated mathematical and statistical technique along with the extensive data collection and functional literature survey enables us to understand various cancers at the fundamental level. The framework considered here focuses on finding important proteins based on their position in the individual networks, which can be extended to those diseases for which very less information is available about the genes which are responsible for the occurrence of the disease. Furthermore, multilayer framework revealing common proteins for different cancers provide a direction for developing novel drugs, therapeutic targets and biomarkers along with the nascent concept of single drug therapy for multiple diseases and personalized medicine in a time efficient and cost effective manner.

## Methods

**Data assimilation and Network construction.** There are two basic components of a network namely, nodes and edges. Here we study PPI networks of the normal and the disease cells where nodes are the proteins and edges denote interactions between the proteins. Nodes in a normal and the corresponding disease network are selected on the basis of their expression in a cell of the normal or disease tissue, respectively. For instance if a protein is expressed as in the normal state of the breast cell, it is considered in the construction of the breast normal network and similarly, if a protein is expressed in the breast tumor (malignant) cell, it is considered in construction of the breast cancer network. After diligent and enormous efforts of mining literature and database text, we collect the list of proteins in the normal tissues and the corresponding cancer tissues from various literature and bioinformatics sources (databanks) namely GenBank<sup>86</sup> and UniProtKB, which mines the proteomic data from various other repositories like European Bioinformatics Institute, the Swiss Institute of Bioinformatics, and the Protein Information Resource etc<sup>87</sup>. We enlist the proteins for a particular tissue by searching relevant keywords, such as its target tissue/cell type in the search panel of various databanks. To keep the authenticity of the data we only take those proteins into account which are reviewed and cited (literature authenticated). Additionally, there are numerous cell lines available for biological studies, but a very few have been exploited for their maximum proteomic insight. We gather the protein expression informations through various cell-lines comprising of different origin (human, mouse, horse etc) from available online literatures to make the data more complete. The details of the cell-lines databanks can be found in the “Data collection and network construction” sub-section of the Supplementary Material. The details of all the proteins for seven different tissues for the normal and disease states can be accessed from<sup>88</sup>. Once all the proteins for seven different tissues for the normal and disease states are collected, leading to fourteen datasets, the interacting partners of these proteins are retrieved from the STRING database version 9.1<sup>89</sup>. An interaction between a pair of proteins is considered if there exists a direct (i.e. physical), indirect (i.e. functional) or both relation between them. STRING provides the physical and functional interactions for a given list of the proteins. Note that in this work, we consider interactions between

a pair of proteins from the STRING database for both the normal and disease networks so any change in the interactions in a particular state (disease or normal) arises only due to deletion or addition of new nodes in the networks. Note that, the information on whether a pair of proteins are really interacting in a particular state, i.e. the dynamical nature of PPIs are missing. In this way, we have seven networks for the normal and seven networks for the corresponding disease states. The protein-protein interactions of all the fourteen networks as adjacency list can be found in ref. 88.

Next, we define the interaction matrix or the adjacency matrix ( $A_{ij}$ ) of the network as,

$$A_{ij} = \begin{cases} 1 & \text{if } i \sim j \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

We investigate the PPI networks for their various structural and spectral properties.

**Properties of Complexome.** *Structural properties.* Several statistical measures have been proposed to understand specific features of a network<sup>90,91</sup>. One of the most basic structural parameter of a network is the degree of a node ( $k_i$ ), which is defined as the number of neighbors of a node has ( $k_i = \sum_j A_{ij}$ ). The degree distribution  $P_k$ , revealing the fraction of vertices with the degree  $k$ , is known to be a fingerprint of the underlying network<sup>90</sup>. Another important parameter of a network is the clustering coefficient (CC) of the network. Clustering coefficient  $CC_i$  of a node (say  $i$ ) can be written as the ratio of the number of interactions the neighbors of a particular node has and number of possible connections the neighbors can have<sup>45,92</sup>. The average of all the individual  $CC_i$  gives the average clustering coefficient ( $\langle CC \rangle$ ). It characterizes the overall tendency of the nodes to form clusters or groups. Further, the betweenness centrality ( $\beta_i$ ) of a node  $i$  is defined as the fraction of shortest paths that pass through the node  $i$ <sup>90</sup>,

$$\beta_i^j = \sum_{st} \frac{n_{st}^i}{g_{st}}, \quad (2)$$

$n_{st}^i$  denotes the number of paths from  $s$  to  $t$  that passes through  $i$  and  $g_{st}$  is the total number of paths from  $s$  to  $t$  in the network. Further, we calculate the diameter of the network which measures the longest of the shortest paths between all the pairs of the nodes<sup>44</sup>.

Another important property of a network that helps us in distinguishing the normal from the disease datasets is degree-degree correlation ( $r$ ). This property measures the tendency of nodes to connect with the nodes having similar number of edges<sup>52,93,94</sup> and can be defined as,

$$r = \frac{\left[ \frac{1}{M} \sum_i j_i k_i \right] - \left[ \frac{1}{M} \sum_i \frac{1}{2} (j_i + k_i)^2 \right]}{\left[ \frac{1}{M} \sum_i (j_i^2 + k_i^2) \right] - \left[ \frac{1}{M} \sum_i \frac{1}{2} (j_i + k_i)^2 \right]}, \quad (3)$$

where  $j_i$  and  $k_i$  are the degrees of the nodes connected through the  $i^{\text{th}}$  edge, and  $M$  is the total number of edges in the network. The value of  $r$  being negative (positive) corresponds to a dis(assortative) network.

Further, to understand the network architecture, we identify the weak ties in the network by calculating the edge-betweenness centrality and overlap of the pair of nodes. Granovetter's weak ties hypothesis: a socially driven network tool highlights the importance of an edge in the network through the strength of their tie by calculating the edge-betweenness centrality ( $\beta_L$ ) with inverse relation to the overlap ( $O$ ) of their neighborhoods<sup>77</sup>. The  $\beta_L$  can be defined as,  $\beta_L = \sum_{v \in V} \sum_{w \in V/v} \sigma_{vw}(e) / \sigma_{vw}$ , where  $\sigma_{vw}(e)$  is the number of shortest paths between  $v$  and  $w$  that contain  $e$ , and  $\sigma_{vw}$  is the total number of shortest paths between  $v$  and  $w$ <sup>79</sup>. Next, the overlap of the neighborhood ( $O_{ij}$ ) of two connected nodes  $i$  and  $j$  is defined as,  $O_{ij} = \frac{n_{ij}}{(k_i - 1) + (k_j - 1) - n_{ij}}$ , where  $n_{ij}$  is the number of neighbors common to both nodes  $i$  and  $j$ <sup>79</sup>. Here,  $k_i$  and  $k_j$  represent the degree of the node  $i$  and  $j$ . Then, we calculate Pearson correlation coefficient ( $O - \beta_L$ ) of  $O_{ij}$  and  $\beta_L$  as,

$$O - \beta_L = \frac{(O_{ij} - \langle O_{ij} \rangle)(\beta_L - \langle \beta_L \rangle)}{\sqrt{(O_{ij} - \langle O_{ij} \rangle)^2} \sqrt{(\beta_L - \langle \beta_L \rangle)^2}} \quad (4)$$

*Spectral properties.* Let us denote eigenvalues of the adjacency matrix by  $\lambda_i$ ,  $i = 1, 2, \dots, N$  such that  $\lambda_1 < \lambda_2 < \lambda_3 < \lambda_1 \dots < \lambda_N$ . Further, in order to understand the evolutionary mechanisms involved in normal and cancer state, that plays an important role in the formation of these PPI networks, we calculate the degenerate eigenvalues in the network. First, we investigate the role of node duplication by identifying the nodes sharing exactly the same neighbors from the corresponding adjacency matrices<sup>61,62</sup>. When (i) two rows (columns) have exactly same entries, it is termed as complete row (column) duplication  $R_1 = R_2$ , (ii) the partial duplication of rows (columns) where  $R_1 = R_2 + R_3$ , where,  $R_i$  denotes  $i^{\text{th}}$  row of the adjacency matrix. The count of zero eigenvalues ( $\lambda_0$ ) provides an exact measure of (i) and (ii) conditions<sup>63</sup>. Further, we calculate degeneracy at minus one eigenvalues ( $\lambda_{-1}$ ) which provides an insight to the complete sub-graphs in the network<sup>95</sup>.

**Multilayer Framework.** Analysis of structural and spectral properties suggests overall similarities between the normal and the disease PPI networks. All the seven disease networks are represented as different layers of a disease multilayer network. Similarly, all the seven normal networks form different layers of a normal multilayer

network leading to the normal multilayer network framework. We extract common nodes from (i) all the normal networks, (ii) all the disease networks and (iii) common between all the disease and all the normal networks (union of (i) and (ii)), and investigate various structural and functional properties of the common proteins referring it as multilayer analysis of these three independent subtractive PPI networks for each cancer (Fig. 2(a)). After extracting common proteins, we find their interacting partners from all the disease datasets and analyze various properties of those proteins which are common in all the disease networks.

**Construction of Erdős Rényi and Configuration networks.** Further, we compare properties of PPI networks with the corresponding ER random networks having the same average degree  $\langle k \rangle$  and size  $N^{90}$ , where nodes are randomly connected with a fixed probability  $p$ , calculated as  $k/N$ . This leads to ER networks having the same density distribution as of the corresponding real networks. We use ER network ensemble of 20 networks for all the properties discussed here.

Additionally, we compare properties of PPI normal and disease networks with the corresponding configuration networks. The configuration model in addition of having the same size and number of connections as of a given network, preserves the degree sequence of the given network, by generating a random network with a given degree sequence of an array of size  $m = \frac{1}{2} \sum_{i=1}^N k_i$ . We construct the configuration model network by taking the degree sequence of various PPI networks as input. Each node of the corresponding configuration model is allotted stubs equal to their degree, and then these stubs are paired with a uniform probability<sup>53,96</sup>. This generates a configuration model for a given degree sequence. We generate 20 such realizations for a given degree sequence.

## References

- Venter, J. C. *et al.* The sequence of the human genome. *Science*. **291**, 1304–1351 (2001).
- Roukos D. H. Genome network medicine: innovation to overcome huge challenges in cancer therapy. *Wiley Interdiscip. Rev. Syst. Biol. Med.* **6**, 201–208 (2014).
- Ramaswamy, S., Ross, K. N., Lander, E. S. & Golub, T. R. A molecular signature of metastasis in primary solid tumors. *Nat. Genet.* **33**, 49–54 (2003).
- Marusyk, A., Almendro, V. & Polyak, K. Intra-tumour heterogeneity: a looking glass for cancer? *Nat. Rev. Cancer*. **12**, 323–334 (2012).
- Garraway, L. A. & Lander, E. S. Lessons from the cancer genome. *Cell*. **153**, 17–37 (2013).
- Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature*. **458**, 719–724 (2009).
- Fisher, R., Puzstai, L. & Swanton, C. Cancer heterogeneity: implications for targeted therapeutics. *Br. J. Cancer*. **108**, 479–485 (2013).
- Burrell, R. A., McGranahan, N., Bartek, J. & Swanton, C. The causes and consequences of genetic heterogeneity in cancer evolution. *Nature*. **501**, 338–345 (2013).
- Portela, A. & Esteller, M. Epigenetic modifications and human disease. *Nat. Biotechnol.* **28**, 1057–1068 (2010).
- Siegel, R., Naishadham, D. & Jemal, A. Cancer statistics, 2012. *CA. Cancer J. Clin.* **62**, 10–29 (2012).
- Torre, L. A. *et al.* Global cancer statistics, 2012. *CA Cancer J Clin.* **65**, 87–108 (2015).
- Dominiotto, M., Tsinoremas, N. & Capobianco, E. Integrative analysis of cancer imaging readouts by networks. *Mol. Oncol.* **9**, 1–16 (2015).
- Creixell, P., Schoof, E. M., Erler, J. T. & Linding, R. Navigating cancer network attractors for tumor-specific therapy. *Nat. Biotechnol.* **30**, 842–848 (2012).
- Futreal, P. *et al.* A census of human cancer genes. *Nat. Rev. Cancer*. **4**, 177–183 (2004).
- Strausberg, R. L., Simpson, A. J. & Wooster, R. Sequence-based cancer genomics: progress, lessons and opportunities. *Nat. Rev. Cancer*. **4**, 409–418 (2003).
- Barabasi, A. L. & Oltvai, Z. N. Network biology: understanding the cell's functional organization. *Nat. Rev. Cancer*. **5**, 101–113 (2004).
- Draghici, S. *et al.* A systems biology approach for pathway level analysis. *Genome Res.* **17**, 1537–1545 (2007).
- Balkwill, F. Cancer and the chemokine network. *Nat. Rev. Cancer*. **4**, 540–550 (2004).
- Wang, Q. *et al.* Community of protein complexes impacts disease association. *Eur. J. Hum. Genet.* **20**, 1162–1167 (2012).
- AlQuraishi, M., Kozytiger, G., Jenney, A., MacBeath, G. & Sorger, P. K. A multiscale statistical mechanical framework integrates biophysical and genomic data to assemble cancer networks. *Nat. Genet.* **46**, 1363–1371 (2014).
- Kar, G., Gurosoy, A. & Keskin, O. Human cancer protein-protein interaction network: a structural perspective. *PLoS Comput. Biol.* **5**, e1000601 (2009).
- Jonsson, P. F. & Bates, P. A. Global topological features of cancer proteins in the human interactome. *Bioinformatics*. **22**, 2291–2297 (2006).
- Creixell, P. *et al.* Pathway and network analysis of cancer genomes. *Nat. Methods*. **12**, 615–621 (2015).
- Barabasi, A. L., Gulbahce, N. & Loscalzo, J. Network medicine: a network-based approach to human disease. *Nat. Rev. Cancer*. **12**, 56–68 (2011).
- Califano, A. Predicting protein networks in cancer. *Nat. Genet.* **46**, 1252–1253 (2014).
- Goh, K. I. *et al.* The human disease network. *Proc. Natl. Acad. Sci. USA* **104**, 8685–8690 (2007).
- Rai, A., Menon, A. V. & Jalan, S. Randomness and preserved patterns in cancer network. *Sci. Rep.* **4**, 6368 (2014).
- Jalan, S. & Bandyopadhyay, J. N. Random matrix analysis of network Laplacians. *Physica A*. **387**, 667–674 (2008).
- Papenbrock, T. & Weidenmüller, H. A. Colloquium: Random matrices and chaos in nuclear spectra. *Rev. Mod. Phys.* **79**, 997–1013 (2007).
- Jalan, S. & Bandyopadhyay, J. N. Random matrix analysis of complex networks. *Phys. Rev. E*. **76**, 046107 (2007).
- Fossion, R., Vargas, G. T. & Vieyra, J. L. Random-matrix spectra as a time series. *Phys. Rev. E*. **88**, 060902 (2013).
- Sarkar, C. & Jalan, S. Social patterns revealed through random matrix theory. *Euro. Phys. L.*, **108**, 48003 (2014).
- Bandyopadhyay, J. N. & Jalan, S. Universality in complex networks: Random matrix analysis. *Phys. Rev. E*. **76**, 026109 (2007).
- Jalan, S., Sarkar, C., Madhusudanan, A. & Dwivedi, S. K. Uncovering randomness and success in society. *PLoS One*. **9**, e88249 (2014).
- Jalan, S. *et al.* Spectral analysis of gene co-expression network of Zebrafish. *Euro. Phys. L.* **99**, 48004 (2012).
- Gibson, S. M. *et al.* Massive-scale gene co-expression network construction and robustness testing using random matrix theory. *PLoS One*. **8**, e55871 (2013).
- Agrawal, A., Sarkar, C., Dwivedi, S. K., Dhasmana, N. & Jalan, S. Quantifying randomness in protein-protein interaction networks of different species: A random matrix approach. *Physica A*. **404**, 359–367 (2014).
- Shinde, P., Yadav, A., Rai, A. & Jalan, S. Dissortativity and duplications in oral cancer. *Eur. Phys. J. B*. **88**, 197 (2015).
- Weinstein, J. N. *et al.* The cancer genome atlas pan-cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).



40. Janga, S. C. & Edupuganti, M. M. R. Systems and Network-Based Approaches for Personalized Medicine. *Curr. Synth. Syst. Biol.* **2**, e109 (2014).
41. Beerenwinkel, N. *et al.* Genetic progression and the waiting time to cancer. *PLoS Comput. Biol.* **3**, e225 (2007).
42. Lorimer, T., Gomez, F. & Stoop, R. Two universal physical principles shape the power-law statistics of real-world networks. *Sci. Rep.* **5**, 12353 (2015).
43. Shinde, P. & Jalan, S. A multilayer protein-protein interaction network analysis of different life stages in *Caenorhabditis elegans*. *Euro. Phys. L.* **112**, 58001 (2015).
44. Albert, R., Jeong, H. & Barabási, A. L. Internet: Diameter of the world-wide web. *Nature*. **401**, 130–131 (1999).
45. Watts, D. J. & Strogatz, S. H. Collective dynamics of small-world networks. *Nature*. **393**, 440–442 (1998).
46. Csirmely, P. & Korcsmros, T. Cancer-related networks: a help to understand, predict and change malignant transformation. *Semin. Cancer Biol.* **23**, 209–212 (2013).
47. Iakoucheva, L. M., Brown, C. J., Lawson, J. D., Obradovi, Z. & Dunker, A. K. Intrinsic disorder in cell-signaling and cancer-associated proteins. *J. Mol. Biol.* **323**, 573–584 (2002).
48. Salido-Guadarrama, I., Romero-Cordoba, S., Peralta-Zaragoza, O., Hidalgo-Miranda, A. & Rodríguez-Dorantes, M. MicroRNAs transported by exosomes in body fluids as mediators of intercellular communication in cancer. *Oncotargets Ther.* **7**, 1327–1338 (2014).
49. Stuart, J. M., Segal, E., Koller, D. & Kim, S. K. A gene-co-expression network for global discovery of conserved genetic modules. *Science*. **302**, 249–255 (2003).
50. Jalan, S., Solymosi, N., Vattay, G. & Li, B. Random matrix analysis of localization properties of gene co-expression network. *Phys. Rev. E*. **81**, 046118 (2010).
51. De Smet, R. & Marchal, K. Advantages and limitations of current network inference methods. *Nat. Rev. Microbiol.* **8**, 717–729 (2010).
52. Jalan, S. & Yadav, A. Assortative and disassortative mixing investigated using the spectra of graphs. *Phys. Rev. E*. **91**, 012813 (2015).
53. Newman, M. E., Strogatz, S. H. & Watts, D. J. Random graphs with arbitrary degree distributions and their applications. *Phys. Rev. E*. **64**, 026118 (2001).
54. Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabási, A. L. Hierarchical organization of modularity in metabolic networks. *Science*. **297**, 1551–1555 (2002).
55. Alon, U. *An Introduction to Systems Biology: Design Principles of Biological Circuits*. CRC press (Chapman & Hall/CRC, London) (2006).
56. Yeger-Lotem, E. *et al.* Network motifs in integrated cellular networks of transcription regulation and protein-protein interaction. *Proc. Natl. Acad. Sci. USA*. **101**, 5934–5939 (2004).
57. Dwivedi, S. K. & Jalan, S. Emergence of clustering: Role of inhibition. *Phys. Rev. E*. **90**, 032803 (2014).
58. Kitano, H. Cancer as a robust system: implications for anticancer therapy. *Nat. Rev. Cancer*. **4**, 227–235 (2004).
59. Jalan, S., Kanhaiya, K., Rai, A., Bandapalli, O. R. & Yadav, A. Network Topologies Decoding Cervical Cancer. *PLoS one*. **10**(8), p.e0135183 (2015).
60. De Aguiar, M. A. M. & Bar-Yam, Y. Spectral analysis and the dynamic response of complex networks. *Phys. Rev. E*. **71**, 016106 (2005).
61. Yadav, A. & Jalan, S. Origin and implications of zero degeneracy in networks spectra. *Chaos*. **25**, 043110 (2015).
62. Kitano, H. Biological robustness. *Nat. Rev. Cancer*. **5**, 826–837 (2004).
63. Golub, G. H. & Van Loan, C. F. *Matrix Computations* (Vol. 3). JHU Press (2012).
64. Nowell, P. C. The clonal evolution of tumor cells. *Science* **194**, 23–28 (1976).
65. Burrell, R. A., McGranahan, N., Bartek, J. & Swanton, C. The causes and consequences of genetic heterogeneity in cancer evolution. *Nature* **501**, 338345 (2013).
66. Hanahan, D. & Weinberg, R. A. The hallmarks of cancer. *Cell*. **100**, 5770 (2000).
67. Yamada, T. & Bork, P. Evolution of biomolecular networks lessons from metabolic and protein interactions. *Nat. Rev. Mol. Cell Biol.* **10**, 791–803 (2009).
68. Typas, A. & Sourjik, V. Bacterial protein networks: properties and functions. *Nat. Rev. Microbiol.* **13**, 559–572 (2015).
69. Allgar, V. L. & Neal, R. D. Delays in the diagnosis of six cancers: analysis of data from the National Survey of NHS Patients: Cancer. *Br. J. Cancer*. **92**, 1959–1970 (2005).
70. Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: the next generation. *Cell*. **144**, 646–674 (2011).
71. Croft, D. *et al.* The Reactome pathway knowledgebase. *Nucleic Acids Res.* **42**, D472–D477 (2014).
72. Batada, N. N., Hurst, L. D. & Tyers, M. Evolutionary and physiological importance of hub proteins. *PLoS Comput. Biol.* **2**, e88 (2006).
73. Okada, H. & Mak, T. W. Pathways of apoptotic and non-apoptotic death in tumour cells. *Nat. Rev. Cancer*. **4**, 592–603 (2004).
74. Fidler, I. J. The pathogenesis of cancer metastasis: the seed and soil hypothesis revisited. *Nat. Rev. Cancer*. **3**, 453–458 (2003).
75. Rak, J. W., St. Croix, B. D. & Kerbel, R. S. Consequences of angiogenesis for tumor progression, metastasis and cancer therapy. *Anticancer drugs*. **6**, 3–18 (1995).
76. Evan, G. I. & Vousden, K. H. Proliferation, cell cycle and apoptosis in cancer. *Nature*. **411**, 342–348 (2001).
77. Granovetter, M. S. The strength of weak ties. *Am. J. Sociol.* **78**, 1360–1380 (1973).
78. Sarkar, C. & Jalan, S. Social patterns revealed through random matrix theory. *Euro. Phys. L.* **108**, 48003 (2014).
79. Onnela, J. P. *et al.* Analysis of a large-scale weighted network of one-to-one human communication. *New J. Phys.* **9**, 179–206 (2007).
80. Xie, Z., Wroblewska, L., Prochazka, L., Weiss, R. & Benenson, Y. Multi-input RNAi-based logic circuit for identification of specific cancer cells. *Science*. **333**, 1307–1311 (2011).
81. Tabassum, D. P. & Polyak, K. Tumorigenesis: it takes a village. *Nat. Rev. Cancer*. **15**, 473–483 (2015).
82. Meltzer, P. S. Cancer genomics: small RNAs with big impacts. *Nature*. **435**, 745–746 (2005).
83. Visone, R. & Croce, C. M. MiRNAs and cancer. *Am. J. Pathol.* **174**, 1131–1138 (2009).
84. Vergoulis, T. *et al.* TarBase 6.0: capturing the exponential growth of miRNA targets with experimental support. *Nucleic Acids Res.* **40**, D222–D229 (2012).
85. Sachdeva, M. & Mo, Y. Y. MicroRNA-145 suppresses cell invasion and metastasis by directly targeting mucin 1. *Cancer Res.* **70**, 378–387 (2010).
86. Benson, D. A. *et al.* GenBank. *Nucleic Acids Res.* **41**, D3642 (2013).
87. Wu, C. H. *et al.* The Universal Protein Resource (UniProt): an expanding universe of protein information. *Nucleic Acids Res.* **34**, D187–D191 (2006).
88. Rai, A. & Jalan, S. Supplementary data: Understanding cancer complexome using networks, spectral graph theory and multilayer framework. figshare. <https://dx.doi.org/10.6084/m9.figshare.4193409.v1> (2016).
89. Franceschini, A. *et al.* STRING v9. 1: protein-protein interaction networks, with increased coverage and integration. *Nucleic Acids Res.* **41**, D808–D815 (2013).
90. Albert, R. & Barabási, A. L. Statistical mechanics of complex networks. *Rev. Mod. Phys.* **74**, 47–97 (2002).
91. Boccaletti, S., Latora, V., Moreno, Y., Chavez, M. & Hwang, D. U. Complex networks: Structure and dynamics. *Phys. Rep.* **424**, 175–308 (2006).
92. Newman, M. E. The structure and function of networks. *Comput. Phys. Commun.* **147**, 40–45 (2002).
93. Newman, M. E. Assortative mixing in networks. *Phys. Rev. L.* **89**, 208701 (2002).



94. Rivera, M. T., Soderstrom, S. B. & Uzzi, B. Dynamics of dyads in social networks: Assortative, relational, and proximity mechanisms. *Annu. Rev. Sociol.* **36**, 91–115 (2010).
95. Van Mieghem, P. *Graph spectra for complex networks*. Cambridge University Press (2010).
96. Molloy, M. & Reed, B. A critical point for random graphs with a given degree sequence. *Random structures & algorithms*. **6**, 161–80 (1995).

### Acknowledgements

S.J. acknowledges Council of Scientific and Industrial Research (CSIR) grant (25(0205)/12/EMR-II) and Department of Science and Technology (DST), Government of India grant EMR/2014/000368 for financial support. AR thanks Amit Kumar Pawar and Sanjiv Kumar Dwivedi for helping with data and codes, respectively and all the members of Complex Systems Lab for useful discussions.

### Author Contributions

S.J. conceived and supervised the project. A.R. collected the data and constructed the networks. S.J., A.R. and P.P. performed the numerical experiments and analyzed the data for various network properties. A.R., J.N. and K.L. did the functional analysis of the important proteins. J.N., K.L. and R.C. performed the miRNA analysis. All the authors wrote and approved the manuscript.

### Additional Information

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing financial interests:** The authors declare no competing financial interests.

**How to cite this article:** Rai, A. *et al.* Understanding cancer complexome using networks, spectral graph theory and multilayer framework. *Sci. Rep.* **7**, 41676; doi: 10.1038/srep41676 (2017).

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

© The Author(s) 2017