

NAPS

Network Analysis of Protein Structures

Reference Manual

December 31, 2018

**Broto Chakrabarty, Varun Naganathan, Kanak Garg, Yash Agarwal and
Nita Parekh,**

Centre for Computational Natural Sciences and Bioinformatics,
International Institute of Information Technology-Hyderabad, India

Contents

1.	About NAPS	3
2.	System requirement	4
3.	Network construction.....	5
4.	Global properties.....	8
5.	Network visualization	10
6.	Node centrality analysis.....	16
7.	Edge centrality:	20
8.	Shortest path analysis.....	21
9.	<i>k</i> -clique analysis.....	22
10.	Graph spectral analysis	23
11.	Multi domain analysis.....	24
12.	Analysis of weighted network	26
13.	Analysis of Protein complex.....	28
14.	RNA Network	31
15.	Protein-Nucleic Acid Complexes	34
16.	Network Analysis of Molecular Dynamics Data	40
17.	References.....	47

1. About NAPS

NAPS is an online portal for the construction and analysis of Protein Contact Network (PCN), also referred as Residue Interaction Network (RIN), or protein structure graph (PSG). A protein structure can be represented as a network by providing the PDB id (or uploading the PDB file) on the portal: <http://bioinf.iiit.ac.in/NAPS/index.php>. The snapshot of the home page is shown in Figure 1.1.

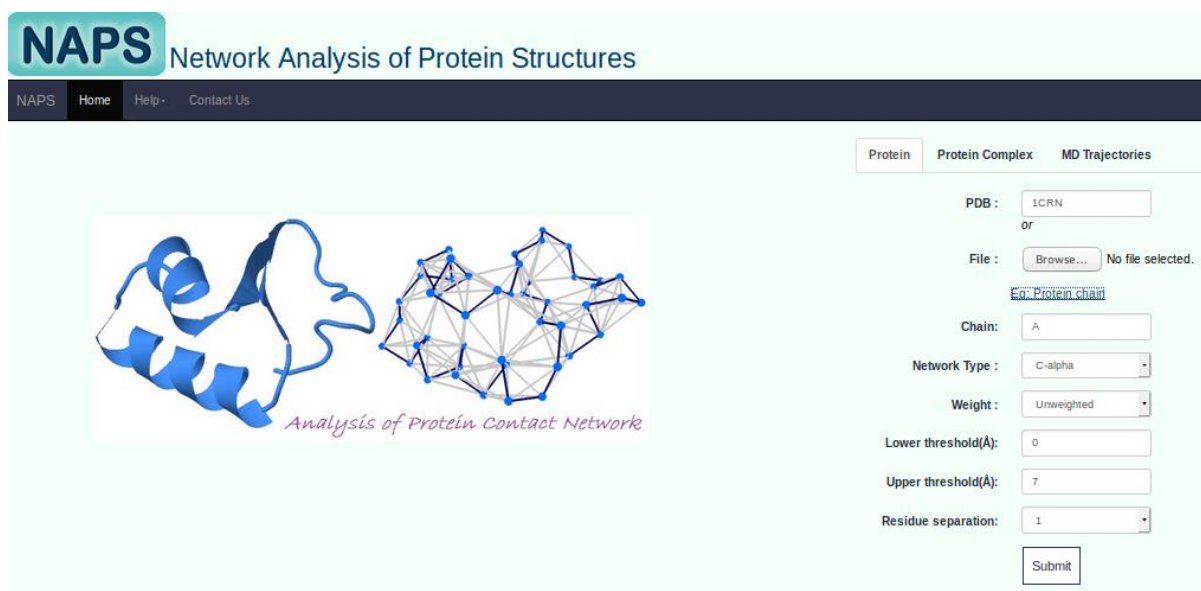


Figure 1.1: Home page of NAPS portal.

2. System requirement

The portal is browser independent and needs no added plugins to be installed in the browser. On providing the PDB id (or upload PDB file), the construction of the network and other computations are performed on the backend server and the data is sent back to the frontend web browsers. The 3D structure of the protein, network visualization and other features are displayed using javascript applets which are browser independent and utilize the client resources (user system resources). Therefore the upper limit of the structure size that can be handled by the portal is determined by the system configuration of the user. For each page, data transfer to/from the server is required only once while loading the page, after which all the activities on the page are performed at the client site with no data transfer required to/from the server. This enhances the performance speed of the portal.

The network visualization is performed using WebGL which is present in all modern browsers. Some of the popular browsers with WebGL are:

Firefox Moxilla: Version 4.0 and above

Google Chrome: Version 9 and above

Safari: Version 6.0 and above installed on OS X Mountain Lion, Mac OS X Lion and Safari 5.1 on Mac OS X Snow Leopard.

The portal displays an error message if WebGL support is not available. If the browser version is above the minimum required version stated above, it implies that WebGL is turned off in the browser. The user needs to manually enable WebGL in such a case.

3. Network construction

The portal provides an option to construct following types of network representations of a protein structure, depending on the type of analysis required.

- a) **C-alpha**: An amino acid residue represented by the C-alpha atom is considered as node in the network and an edge is constructed if the distance between a pair of C-alpha atoms is within the lower and upper thresholds defined by the user (default upper threshold = 7 Å; lower threshold = 0 Å).

- i. **Unweighted**: All edges are considered equally important.
- ii. **Weighted**: Edge weight for a C-alpha weighted network is given by:

$$w_{ij} = \frac{1}{d_{ij}}$$

where d_{ij} is the euclidean distance between C-alpha atoms of i^{th} and j^{th} residues.

- b) **C-beta**: An amino acid residue represented by the C-beta atom is considered as node in the network and an edge is constructed if the distance between the C-beta atoms (C-alpha for GLY) is within the lower and upper thresholds defined by the user (default upper threshold = 7 Å; lower threshold = 0 Å).

- i. **Unweighted**: All edges are considered equally important.
- ii. **Weighted**: Edge weight weighted network is given by:

$$w_{ij} = \frac{1}{d_{ij}}$$

where d_{ij} is the euclidean distance between C-beta atoms of i^{th} and j^{th} residues.

- c) **Atom pair contact**: An amino acid residue is considered as node in the network and an edge is constructed if the distance between any pair of atoms of the residue pair is within the lower and upper thresholds defined by the user (default upper threshold = 5 Å; lower threshold = 0 Å).

- i. **Unweighted**: All edges are considered equally important.
- ii. **Weighted**: Edge weight is given by the number of atom pairs within cutoff distance.

- d) **Centroid (centre of mass)**: An amino acid residue is considered as node in the network and an edge is constructed if the distance between centre of mass of the residue pair is within the lower and upper thresholds defined by the user (default upper threshold = 8.5 Å; lower threshold = 0 Å).

- i. **Unweighted**: All edges are considered equally important.
- ii. **Weighted**: Edge weight weighted network is given by:

$$w_{ij} = \frac{1}{d_{ij}}$$

where d_{ij} is the euclidean distance between centre of mass of i^{th} and j^{th} residues.

- e) **Interaction strength**: An amino acid residue is considered as node in the network and an edge is constructed if the interaction strength between two residues is more than the threshold defined by the user (default = 4%) (1). The interaction strength as proposed by Brinda and Vishveshwara (2) is calculated as:

$$I_{ij} = \frac{n_{ij}}{\sqrt{N_i * N_j}} * 100$$

where, n_{ij} is the number of side chain atom pairs of the residues i and j within 4.5 Å. N_i and N_j are the normalization values of the residues i and j given by Kannan and Vishveshwara (1) as shown in Table 3.1.

- i. **Unweighted:** All edges are considered equally important.
- ii. **Weighted:** The interaction strength (I_{ij}) is considered as edge weight.

The choice of network type and threshold depend on the biological problem to be addressed. The different network types and some example problems for which they have been used in the past, are listed in Table 3.2. The analysis of one protein structure as an unweighted network is shown in sections 4 – 11. The advanced options for analysis of protein complex are provided in section 12.

Long range interaction network (LIN)

A LIN is constructed by considering edges between residues that are sequentially separated by 10-12 residues along the protein backbone. LIN can be constructed by selecting a suitable threshold residue separation at the home page. A range of 1 to 15 is available as threshold for residue separation. Minimum residue separation of 1 is used as default, where edge is drawn between any pair of residues satisfying the criteria of network construction, including the adjacent residues of the protein backbone. Threshold of 2 or 3 is usually used to remove noise in the network.

Table 3.1: Normalization value for amino acid residues used to construct Interaction Strength Network.

Residue Type	Norm (N)
Alanine	55.7551
Arginine	93.7891
Asparagine	73.4097
Aspartic acid	75.1507
Cystine	54.9528
Glutamine	78.1301
Glutamic acid	18.8288
Glycine	47.3129
Histidine	83.7357
Isoleucine	67.9452
Leucine	72.2517
Lysine	69.6096
Methionine	69.2569
Phenyl alanine	93.3082

Proline	51.3310
Serine	61.3946
Threonine	63.7075
Trptophan	106.703
Tyrosine	100.719
Valine	62.3673

Table 3.2: Different network types and the purpose they are used.

Network Type	Edge weight	Purpose
C_{α}	Unweighted	Analysis of global network properties (3, 4), Protein folding kinetics (5), Analysis of inter- and intra-molecular, communications (6), Identification of proteins with similar folds (7), Structural repeat identification (8–10)
C_{β}	Unweighted	Protein dynamics (11)
Atom pair contact	Unweighted	Protein fold (12)
	Weighted	Network analysis based on physicochemical properties (13)
Centroid	Unweighted	Protein core and exposed residue analysis (14)
Interaction strength	Weighted	Structural stability (2), Identifying side chain clusters (1)

4. Global properties

The global properties of the network are displayed in the network property page along with details of network construction. This is the first page shown after construction of the network. The page can be browsed by selecting 'Properties' option from the pull-down menu of 'Analysis' tab at the top menu bar. The snapshot of the page for an example protein 1CRN (chain A) is shown in Figure 4.1.

The following global parameters are displayed on the network property page:

- a) **N_r**: Number of nodes in the network. This represents the number of residues in the protein (or complex).
- b) **N_e**: Number of edges in the network.
- c) **D**: Shortest path distance between the pair of farthest nodes in the network.
- d) **R**: Shortest path distance of the centre(s) of the network to the farthest node.
- e) **k**: Average degree of all the nodes in the network.
- f) **l**: Average of shortest path between all node pairs in the network.
- g) **c**: Clustering coefficient of the network

The following network analysis links are available on the network property page:

- a) View Network: Network visualization as described in section 5.
- b) Node centrality: Centrality analysis as described in section 6.
- c) Edge centrality: Edge centrality as described in section 7.
- d) Shortest path: Shortest path analysis as described in section 8.
- e) *k*-clique: *k*-clique analysis as described in section 9.
- f) Graph spectra: Graph spectral analysis of adjacency and laplacian matrices as described in section 10.
- g) Domain view: Visual analysis of multi-domain proteins as described in section 11.

The user can also select the type of analysis to be performed from the pull-down menu of the tab 'Analysis' from menu bar.

Note: The global parameters, node centrality and edge centrality are not available if the network has more than one connected components. A network is called one connected component, if at least one path exists between each node pair of the network.

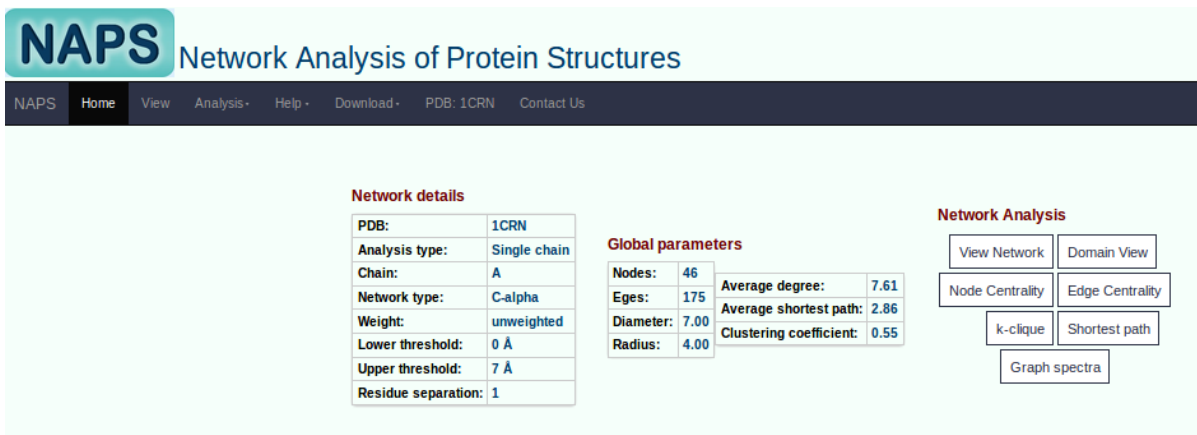


Figure 4.1: Properties of the Network.

The download option on the menu bar of this page provides an option to download:

- a) Global parameters in a text file.
- b) Edgelist file.

5. Network visualization

The network view page is dedicated to different options for visualizing the network along with the three dimensional structure of the protein. The network view page can be accessed from other analysis pages by selecting ‘View’ option from the top menu bar. The 3D structure of the protein is shown using open source javascript based applet, JSmol (15). There are four types of visualization options available in NAPS:

a) 3D structural view

The 3D structure of the protein is shown using open source javascript based applet, JSmol (15).

b) Network View

Depending on the type of network construction chosen, a 3D graphical view is shown in the left panel. To correlate the network to the protein structure, the nodes are plotted using the actual coordinates of the representative atoms (C-alpha/C-beta/centroid) of the corresponding residues in the PDB file. The view page with network view and 3D structure view of the protein 1CRN (chain A) is shown in Figure 5.1.

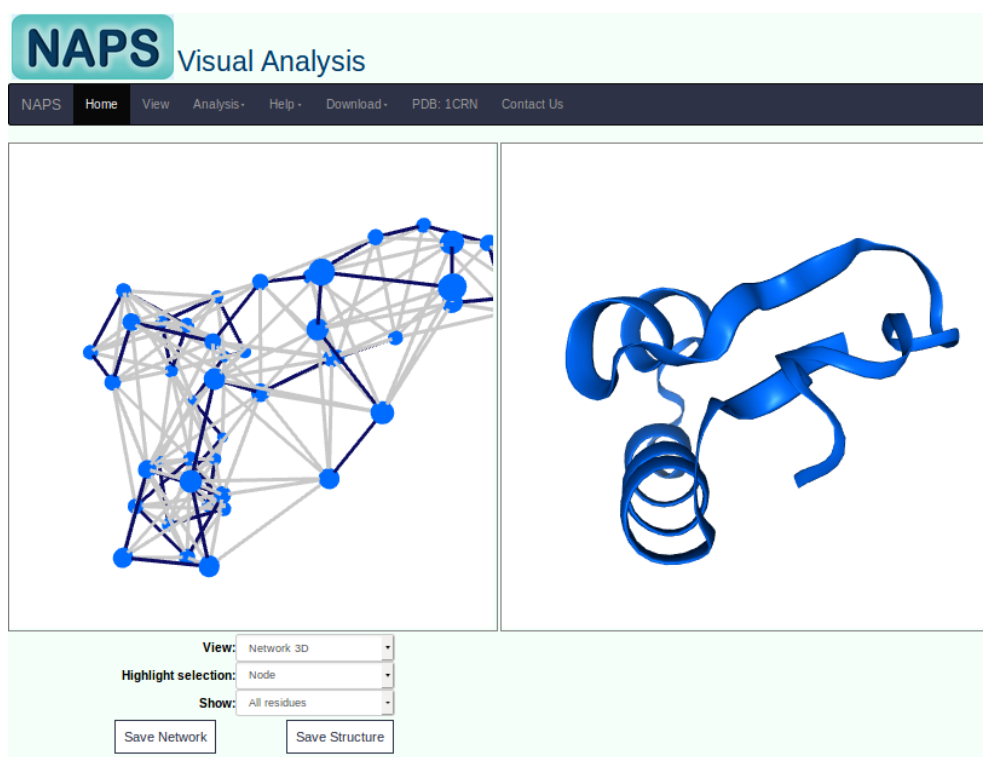


Figure 5.1: Network view page with 3D network representation.

The nodes in the network view are represented in blue color, while two types of edges are shown, dark blue and grey. The dark blue edges trace the backbone connectivity of the protein structure, while the grey edges represent all other nodes within a user defined cutoff threshold of a node. On hovering the mouse over a node in the network view, the node is highlighted with pink color and the chain number, residue number and residue type corresponding to the node are displayed, as shown in Figure 5.1. Two highlight selection options are provided on the network view page:

- i. **Highlight node:** On selecting a node(s) by clicking, it gets highlighted with red color and the corresponding residue on the 3D structure representation also gets highlighted with red color as shown in Figure 5.2.

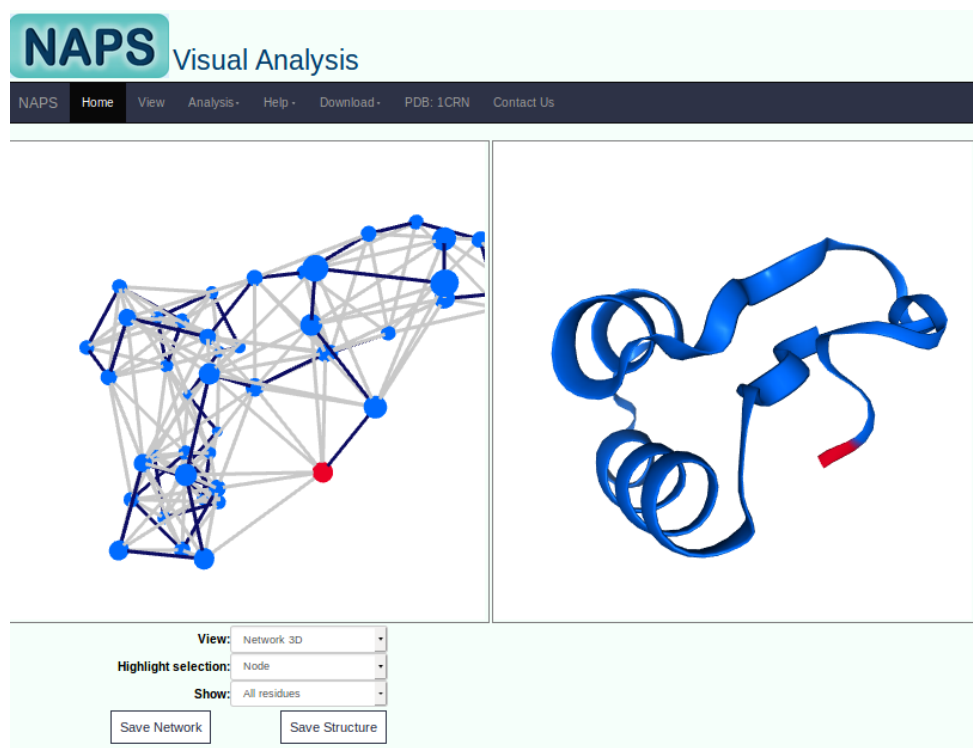


Figure 5.2: 3D representation with node selection.

- ii. **Highlight neighbor:** In a dense network like PCN with many grey edges originating from a node, it is desirable to have easy identification of the immediate neighbors of a node. On selecting a node with highlight selection as neighbor, the node and the corresponding residue in JSmol applet are highlighted in red color, and its immediate neighboring nodes in the network view and the corresponding residues in JSmol applet are highlighted in yellow as shown in Figure 5.3. The residue ids of the selected node and its neighbors are displayed below the JSmol applet.

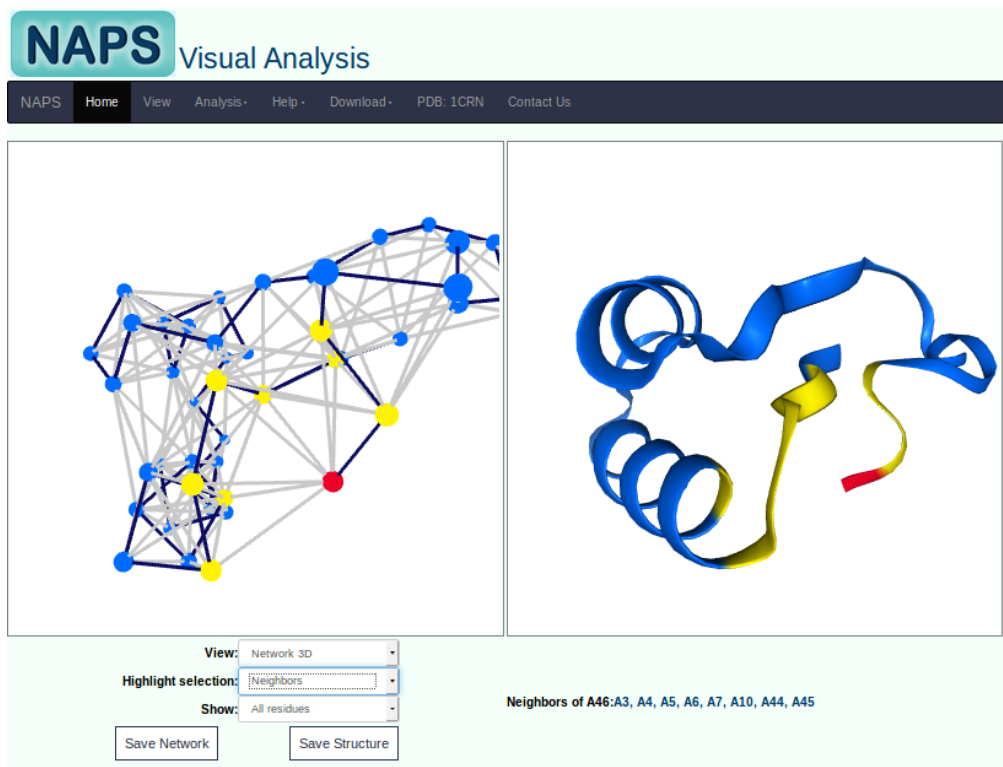


Figure 5.3: 3D Network view with highlight selection as neighbors.

c) Contact map view

Contact map view is a two-dimensional dot matrix representation of the network where a dot (i, j) represents an edge between the i^{th} and j^{th} nodes and is shown in Figure 5.4. The contact map view helps in visual inspection of interaction patterns and the long range interactions within the protein structure. The pattern of connectivity within a secondary structure remains conserved which can be observed in Figure 5.4 where the edges within helix are highlighted in red color. The residue id, chain and residue type of the two amino acids forming the edge are displayed on taking the mouse on an edge in the contact map view. On selecting the edge, the corresponding residues forming the edge are highlighted by red as shown in Figure 5.4. An option to add grid lines at intervals of 10 residues is provided as checkbox.

The download option on the menu bar of the visualization page provides allows download of:

- a) Network view image in PNG format.
- b) Structure view image in PNG format –Network view or Contact map, whichever is selected.
- c) Edgelist - a text file listing the node pairs sharing an edge between them.

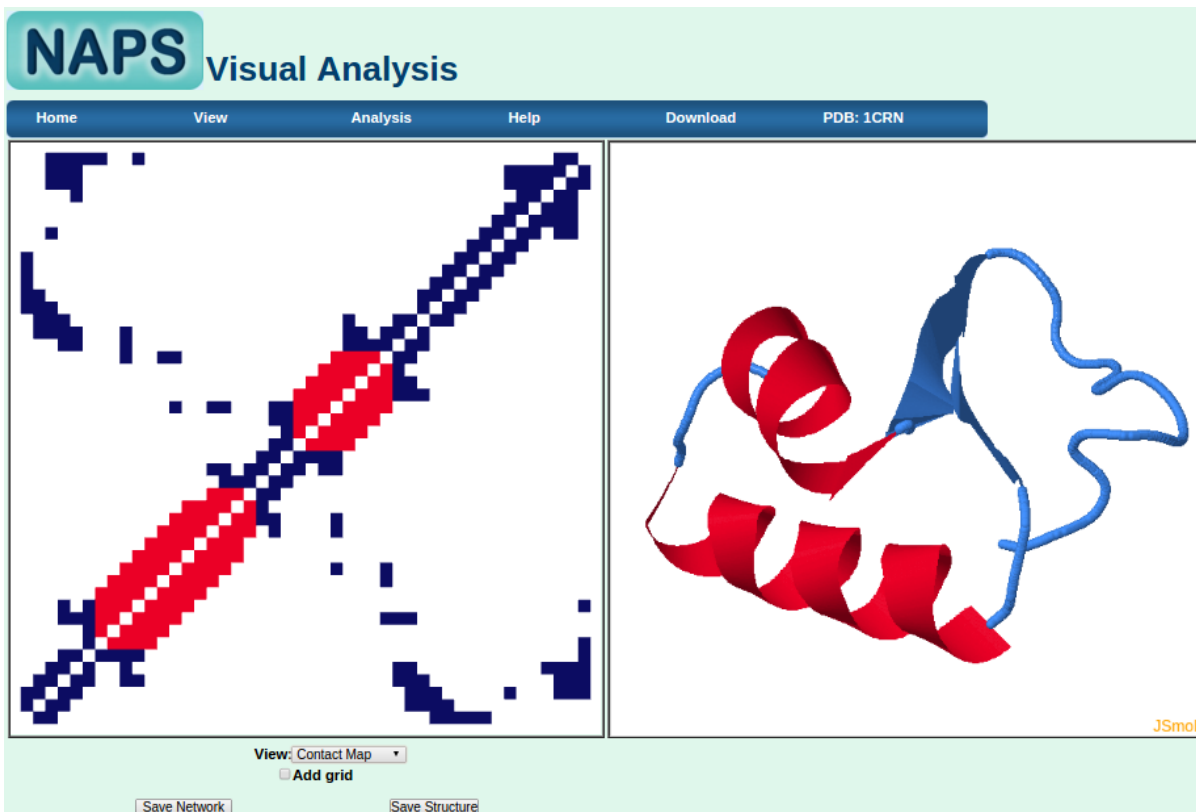


Figure 5.4: Contact Map view.

d) Distance matrix view

The distance matrix captures the distance between the amino acid residues. The element (i,j) of the matrix represents the distance (in Å) between i^{th} and j^{th} residues. The distance matrix view provides a 2D representation of the distance matrix with color representing the distance between the residues as shown in Figure 5.5. The residue id, chain and residue type of the two amino acids represented by a cell in the distance matrix, are displayed on taking the mouse on the distance matrix cell. An option to add grid lines at intervals of 10 residues is provided as checkbox.

Note: Contact map and distance matrix views are available for systems with up to 750 residues.

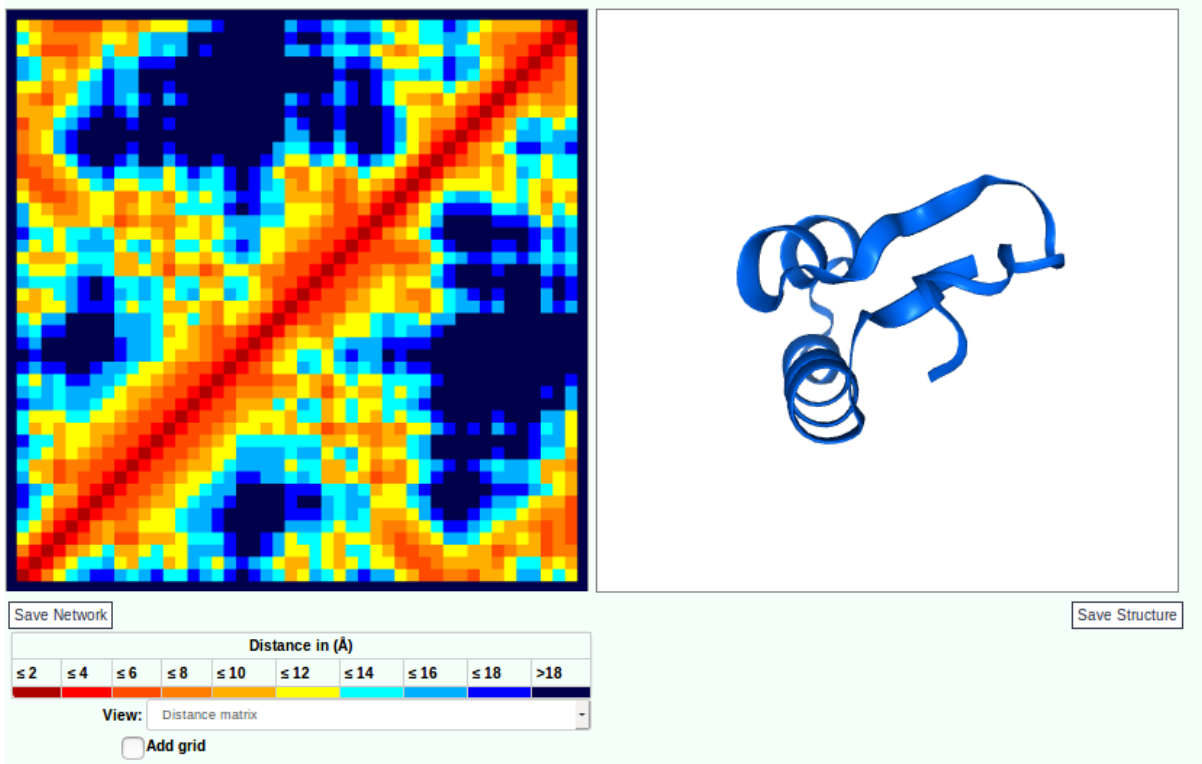


Figure 5.5: Distance matrix representation of the network.

Sub-network view

Visual analysis of a sub-network based on physicochemical properties of the residues can be performed. The network 3D view provides an option to select specific residue types: hydrophobic, hydrophilic and charged. On selecting one of the residue type, all the residues with the selected physicochemical property are highlighted in the network 3D view and the JSmol applet, and edges are shown only between the highlighted residues, as shown in Figure 5.6 for hydrophobic residues.

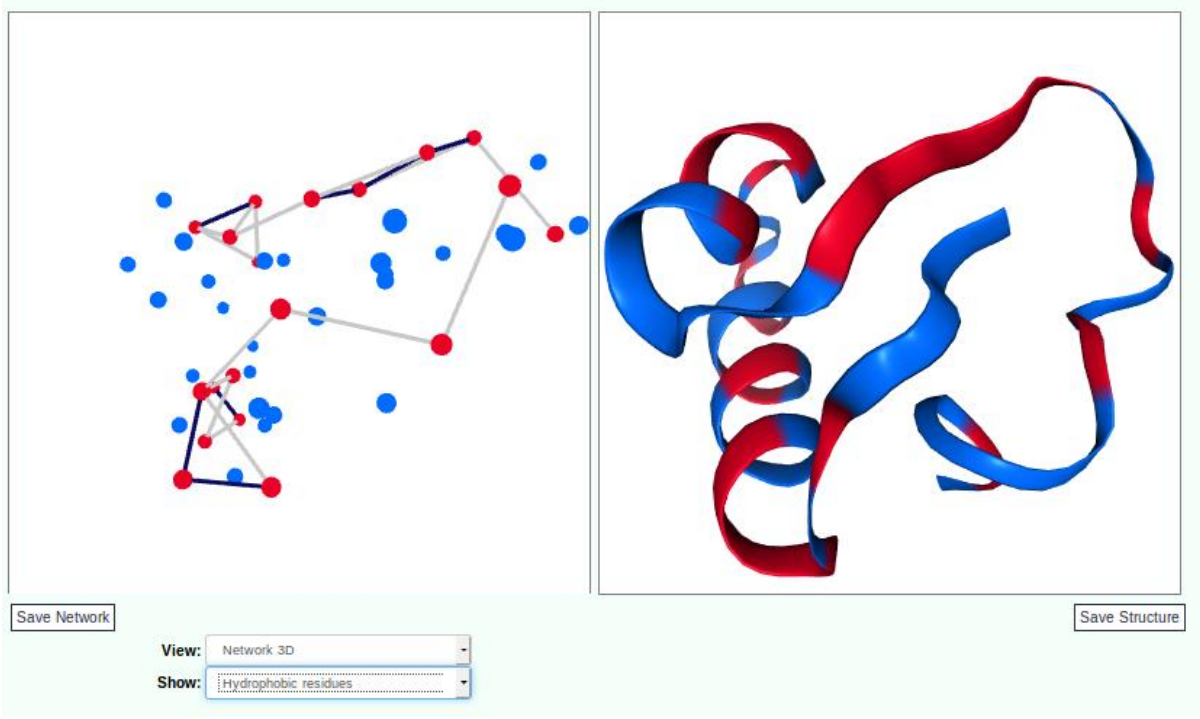


Figure 5.6: Sub-network representation of network 3D view showing Hydrophobic residues and edges between them.

6. Node centrality analysis

Centrality measure of a node provides a quantification of the topological importance of the node in the network. Different centrality measures have been proposed for ranking the nodes in a complex network and quantifying their relative importance. The analysis of centrality measures for a PCN can be performed by selecting the ‘Node centrality’ option from the pull-down menu of ‘Analysis’ tab at the top menu bar. Here we provide option to compute seven node-based centrality measures:

- a) **Degree:** Number of direct neighbors of a node.
- b) **Closeness:** It is the inverse of total shortest path distance of the node to all other nodes of the network.
- c) **Betweenness:** It is the ratio of shortest paths passing through the node.
- d) **Clustering coefficient:** It is the ratio of number of connected neighbors to the total number of connections possible between the neighbors.
- e) **Eigenvector centrality:** It is the eigenvector component corresponding the largest eigenvalue of the adjacency matrix.
- f) **Eccentricity:** Shortest path distance of the node to the farthest node in the network.
- g) **ANN degree:** Average of degree of its immediate neighbors.
- h) **Strength:** Weighted degree represented by cumulative weights of all the edges connected to a node. This is applicable only for weighted networks.

On the centrality analysis page, a table lists the centrality value for each residue for a chosen centrality measure, on the right panel. The interactive network and 3D structure view are displayed on the left panel. The user can select the centrality measure from a drop down menu, can sort it based on the residue id or in descending order of the centrality value. One can also analyze the centrality value of any node by selecting it on the interactive network. The corresponding residue in JSmol applet is simultaneously highlighted along with the corresponding entry in the table, in red color. Snapshot of the interactive analysis of centrality is shown in Figure 6.1.

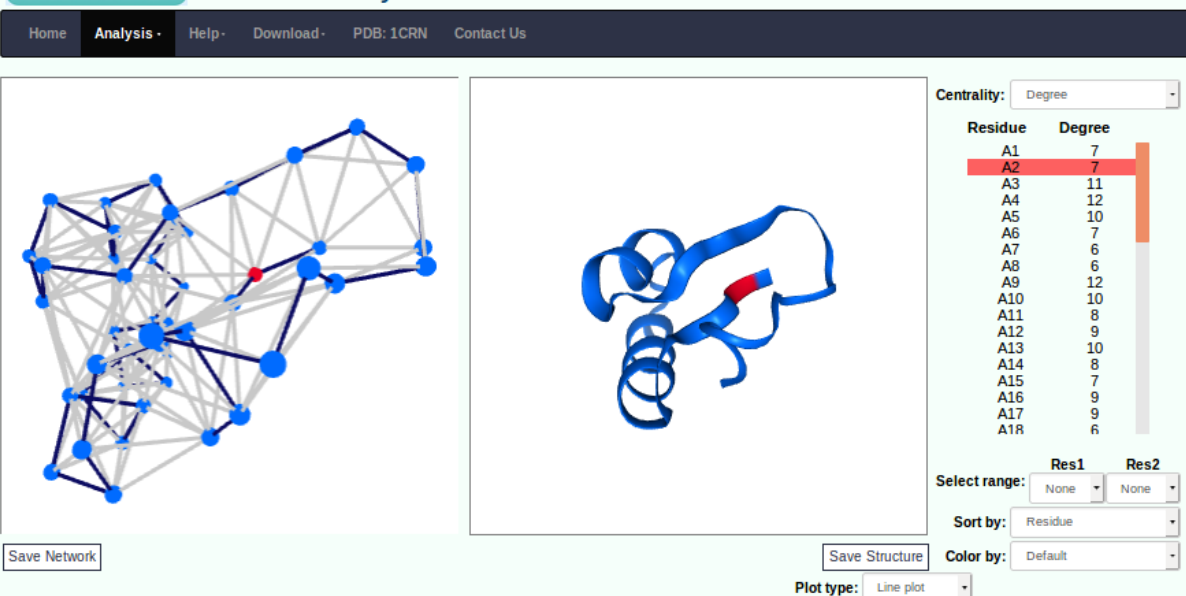


Figure 6.1: Centrality page showing degree centrality with residue A2 highlighted.

For each centrality measure, the following options are provided to aid in the analysis of centrality measures:

- a) **Color by centrality:** By changing the color option from default to centrality, the nodes in the network view and the corresponding residues in the JSmol applet are colored according to the centrality values in a gradient of red (maximum) to yellow (minimum), as shown in Figure 6.2. This helps in easy identification of residues with high centrality values, which are likely to be the residues important for the 3-dimensional fold of the protein, or functionally important (16).

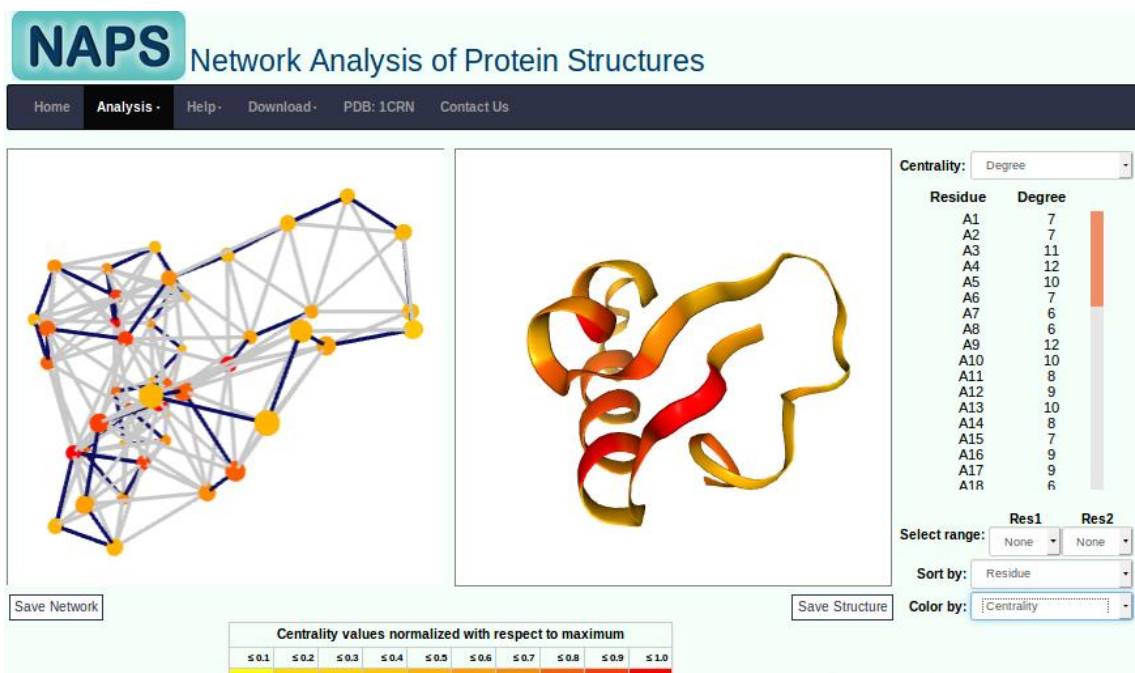


Figure 6.2: Nodes colored according to centrality values.

- b) **Color by Hydrophobicity:** By changing the ‘Color by’ option to ‘Hydrophobicity’, the nodes in the network view and the corresponding residues in the JSmol applet are colored according to the hydrophobicity values of the amino acid residue. The hydrophobicity indices for 20 amino acids given by Kyte and Doolittle range between -4.5 to 4.5 (17). These values are colored in gradients of red to show hydrophobic residues and gradients of blue to show hydrophilic residues, as shown in Figure 6.3.

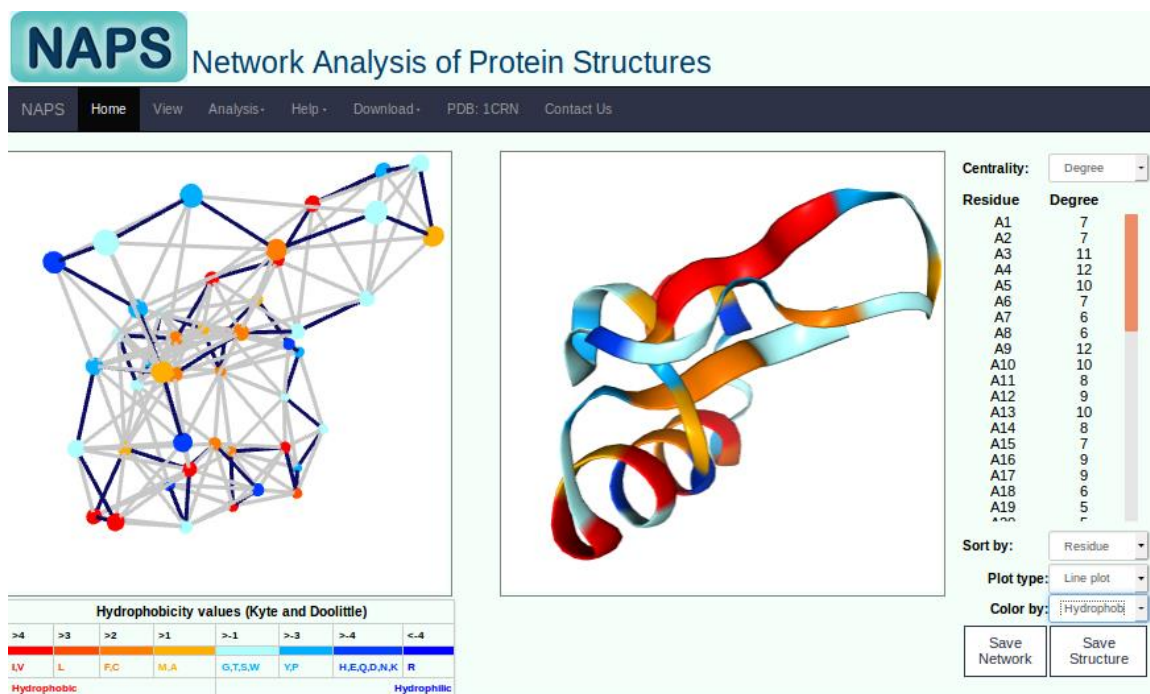


Figure 6.3: Nodes colored according to hydrophobicity.

- c) **Centrality Plots:** By clicking the ‘Generate plot’ button, one can obtain a plot of the centrality measure along the length of the protein, as shown in Figure 6.4. A new page opens where the user can select the file type, line color and resolution of the plot. The user can generate high quality images of the centrality plot for publication purposes and save it by clicking on the download button.
- d) **Highlight residue range:** A set of specific residues or a range of residues can be highlighted by selecting the residue ids or the range from the drop down menu named ‘Select range’ provided with the ‘Color by’ option as ‘default’. The high centrality nodes obtained from the sorted centrality table can be highlighted using this option.
- e) **Highlight neighbor:** It is desirable to identify the immediate neighbors of high centrality nodes of a network. Clicking a node with highlight selection as ‘neighbor’, the node is highlighted by red and its neighbors are highlighted by yellow color. The feature is described in detail in Network Visualization section.

From the download option on the menu bar of this page the user can download:

- Centrality values in tab separated text file.
- Plots of Centrality measures in PNG format.
- Network image in PNG format.

- d) Structure image in PNG format.
- e) Edgelist in tab separated text file.



Figure 6.4: Plot of degree centrality.

7. Edge centrality:

The analysis of edge centrality measures for a PCN can be performed by selecting the 'Edge centrality' option from the pull-down menu of 'Analysis' tab at the top menu bar. The contact map view of the network is displayed along with the protein 3D structure and a table showing the edge betweenness values. On selecting an edge on the contact map, the corresponding edge in the contact map, the two residues in the JSmol applet and the edge betweenness value in the table get highlighted by red color as shown in Figure 7.1.

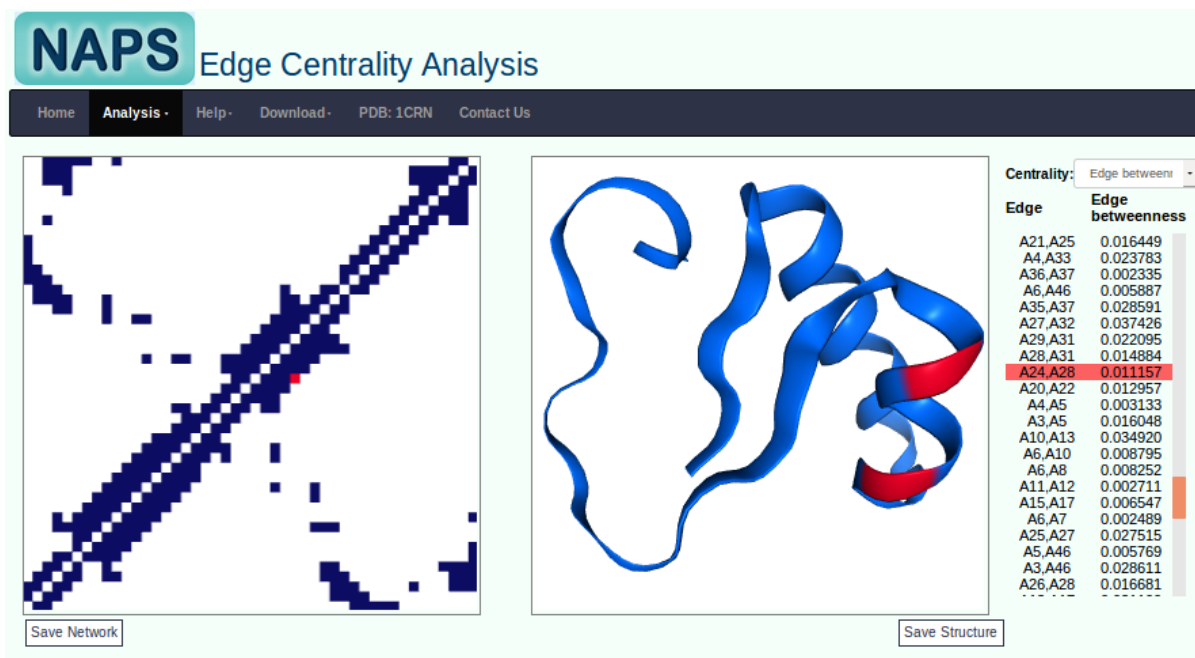


Figure 7.1: Edge centrality page showing edge betweenness.

8. Shortest path analysis

Shortest path distance between two nodes in a network is the minimum number of nodes that need to be traversed in order to reach from one node to the other. The shortest path analysis can be performed by selecting the ‘Shortest path’ option from the pull-down menu of ‘Analysis’ tab at the top menu bar. The user can select two residues from the drop down list, all shortest paths (if more than one) between the two nodes are identified and listed with radio buttons as shown in Figure 8.1. Selecting any path by clicking a radio button highlights the nodes and edges along that path in orange color in both the network view and JSmol applet (Figure 8.1).

The interface displays a network graph on the left and a 3D ribbon structure on the right. The network graph shows nodes as blue spheres and edges as grey lines, with a specific path highlighted in orange. The 3D structure shows a blue ribbon with the same path highlighted in orange. Below the graphs, there is a control panel with two dropdown menus for 'Residue 1' (set to A27) and 'Residue 2' (set to A46), a 'Compute shortest path' button, and 'Save Network' and 'Save Structure' buttons. To the right of the control panel, the 'Shortest path length: 3' is displayed, followed by a list of paths with radio buttons. The first path, 'A27-A33-A3-A46', is selected with a checked radio button.

Figure 8.1: Shortest path analysis page highlighting one of the paths between residue A27 and A46.

From the download option on the menu bar of this page the user can download:

- Shortest path between selected nodes in text format.
- Network image in PNG format.
- Structure image in PNG format.
- Edgelist file.

9. k -clique analysis

A k -clique is a sub-network of k nodes with all the k nodes are connected to each other. These nodes may have edges to nodes outside the sub-network. The clique analysis can be performed by selecting the ‘ k -clique’ option from the pull-down menu of ‘Analysis’ tab at the top menu bar. On selecting k , from the drop down menu, all the cliques of size k are displayed as radio buttons. Selecting a radio button, the nodes and edges of that k -clique are highlighted in orange color in both the network view and JSmol applet as shown in Figure 9.1.

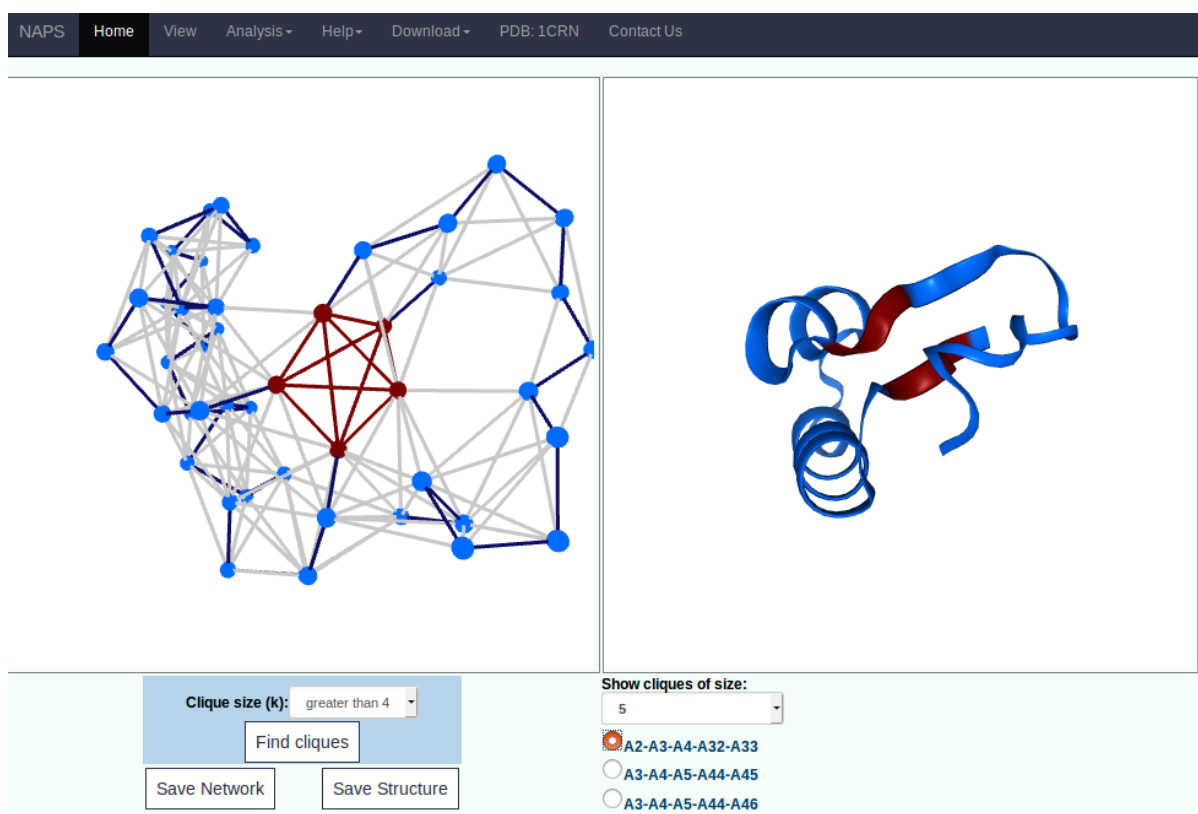


Figure 9.1: Clique of size 5 shown.

From the download option on the menu bar of this page the user can download:

- All clique of selected size in text format.
- Network image in PNG format.
- Structure image in PNG format.

10. Graph spectral analysis

The graph spectral analysis provides spectral analysis of adjacency and laplacian matrices. It can be performed by selecting the 'Graph spectra' option from the pull-down menu of 'Analysis' tab at the top menu bar. The eigenvector component corresponding the largest eigen value of the adjacency matrix and the eigenvector corresponding second smallest eigen value of the laplacian matrix can be analyzed. All the view options discussed for node centrality analysis are available for graph spectral analysis. An example case with nodes in the network view and residues in JSmol colored according to the eigenvector corresponding second smallest eigen value of the laplacian matrix are shown in Figure 10.1.

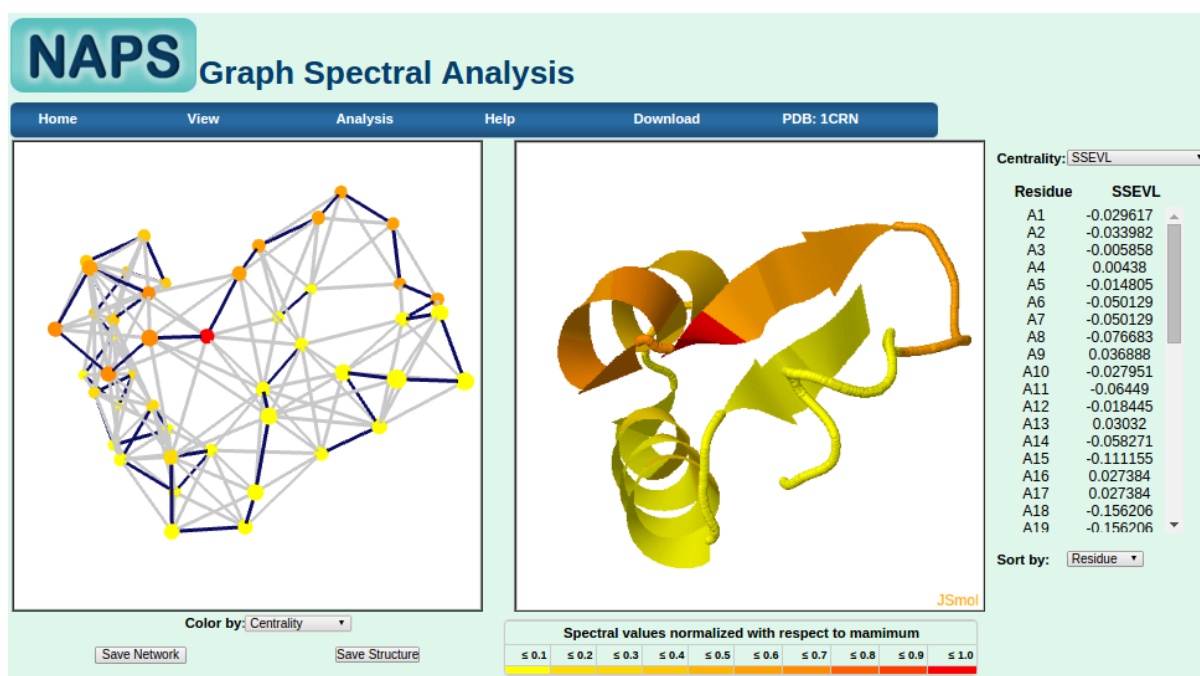


Figure 10.1: Nodes in network view and residues in JSmol applet colored according to the eigenvector corresponding second smallest eigen value of the laplacian matrix.

11. Multi domain analysis

In large multi-domain proteins, it is desirable to analyze the interaction patterns within and between the domains, which is provided in the 'Domain view' page of NAPS. Multi domain analysis can be performed by selecting the 'Domain view' option from the pull-down menu of 'Analysis' tab at the top menu bar. NAPS provides options to select up to 5 domains together. The user can provide the desired number of domains along with the coordinates, which will be used to color the nodes in the network 3D view and the residues in JSmol view. The analysis can be carried out for both contiguous and non-contiguous domains, which is one of the most useful feature required for visual analysis of multi-domain proteins. An example protein, 16PK with two domains is shown in Figure 11.1.

The coordinates of a contiguous domain can be given by first selecting the chain from the drop down menu and then providing the range. For example, input range '5:192' means residues 5 to 192. The coordinates of a non-contiguous domain can be given by providing the ranges separated by ','. For example, '1:10,25,30:35' means residues 1 to 10, residue 25 and residues 30 to 35.

The screenshot displays the NAPS Domain View interface for protein 16PK. The top navigation bar includes 'Home', 'Analysis -', 'Help -', 'Download -', 'PDB: 16PK', and 'Contact Us'. The main content area is split into two panels: a network view on the left and a ribbon structure view on the right. The network view shows a complex graph of nodes and edges, with two distinct clusters of nodes. The ribbon structure view shows the protein structure with two domains highlighted in different colors: cyan and yellow. Below the views are buttons for 'Save Network' and 'Save Structure'. A configuration panel at the bottom allows for domain selection, showing 'Number of domains: 2', 'Domain 1: A 5:192', and 'Domain 2: A 199:406'. An example range '1:10,25,30:35' and a 'Color domains' button are also visible.

Figure 11.1: Domain view showing two domains of the protein 16PK.

The betweenness centrality analysis of the two domain protein is shown in Figure 11.2. It can be observed that the residues in the interface of the two domains have high betweenness values indicating their importance in the inter-domain communications.

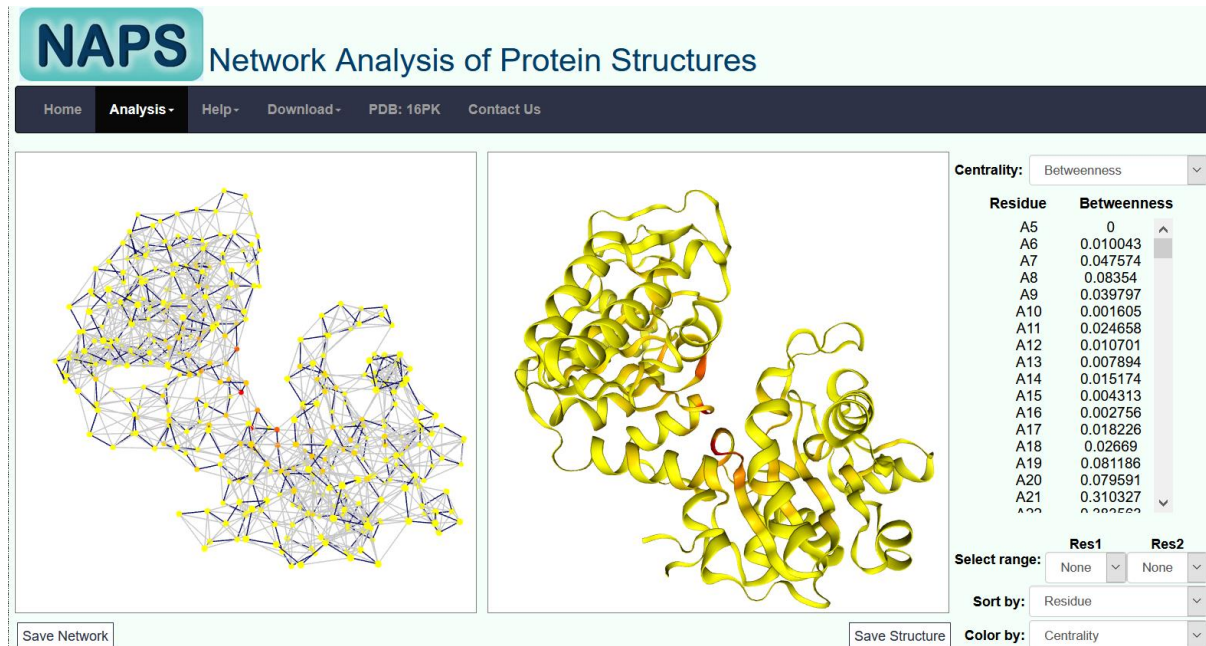


Figure 11.2: Betweenness centrality analysis of multi-domain protein, 16PK.

12. Analysis of weighted network

A weighted network can be constructed by selecting 'weight' option on the home page as 'weighted'. The definition of edge weight for all 5 network types are described in Section 3 (Network construction). All the analysis for a weighted network can be performed similar to the unweighted network as discussed in the previous sections. In this section, we show the analyses pages which are different for a weighted network. The edge lines in 3D network view and the cells in contact map view are color coded according to the edge weights, as shown in Figure 12.1 and Figure 12.2.

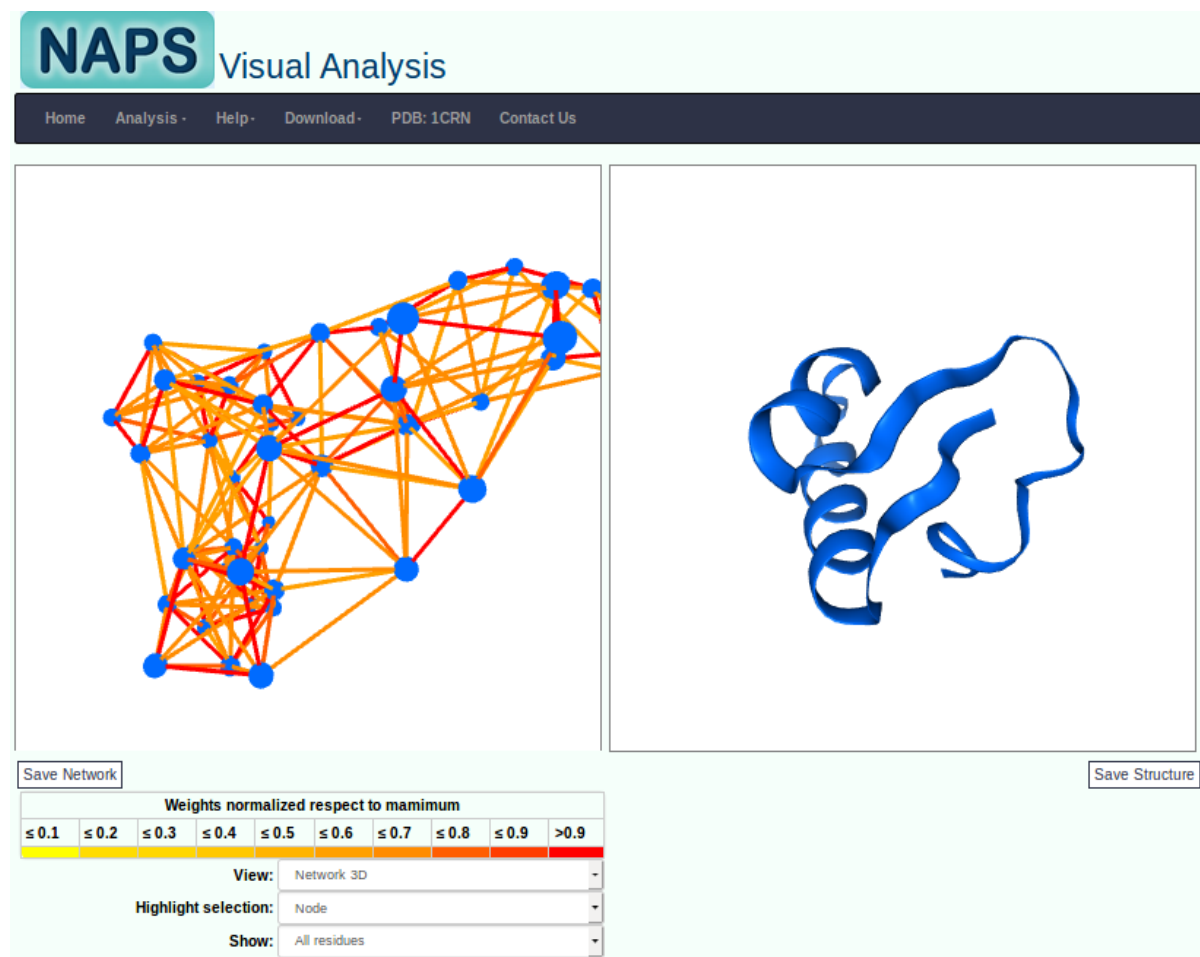


Figure 12.1: Network 3D view showing edge color based on weights.

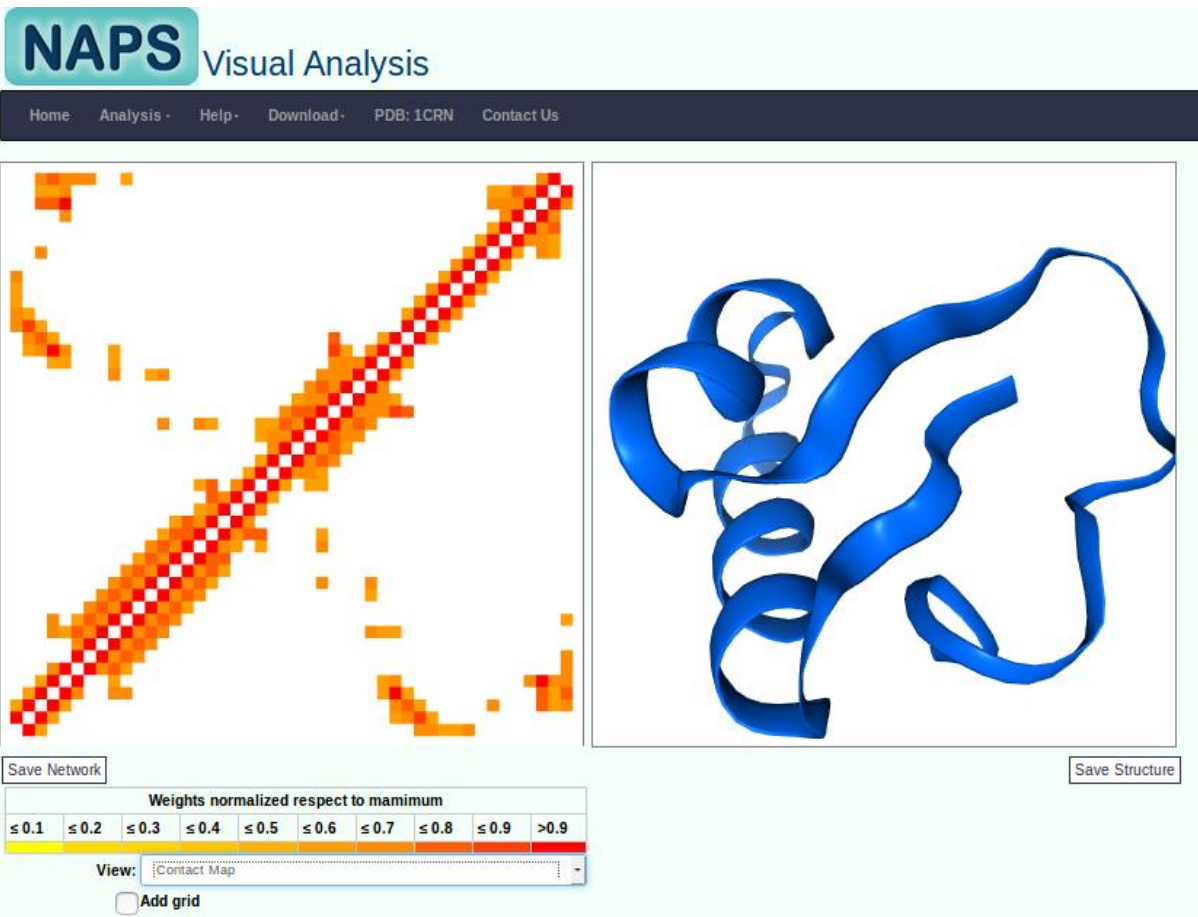


Figure 12.2: Contact map view with cell color showing the edge weights.

Shortest path analysis in weighted network

The inferences from shortest path analysis are based on the principle that the most important path (route) between a pair of nodes is the one with the minimum distance, i.e. the edges with lower distances (weights) are more important. In a weighted Protein Contact Network (Section 3), the weights are defined in such a way that the edges with more weights are more important. In order to make meaningful inferences from shortest path analysis of weighted PCNs, the distance between a pair of nodes i and j is taken as $(1/w_{ij})$ for all shortest path analyses: average shortest path, closeness centrality, betweenness centrality and shortest path analysis between a pair of nodes.

13. Analysis of Protein complex

All the analysis that can be performed on one protein can be performed on a protein complex with up to four chains representing either dimers/trimers/tetramers of the same protein, or a complex of 2-4 interacting proteins. The 3D network view shows the nodes corresponding the four protein chains in four colors with backbone edges in the corresponding colors. Inter-chain edges, i.e. the edges between nodes of different chains are represented in yellow color while all intra-chain edges are colored grey. The snapshots of the tool for different analysis pages are shown in following figures.

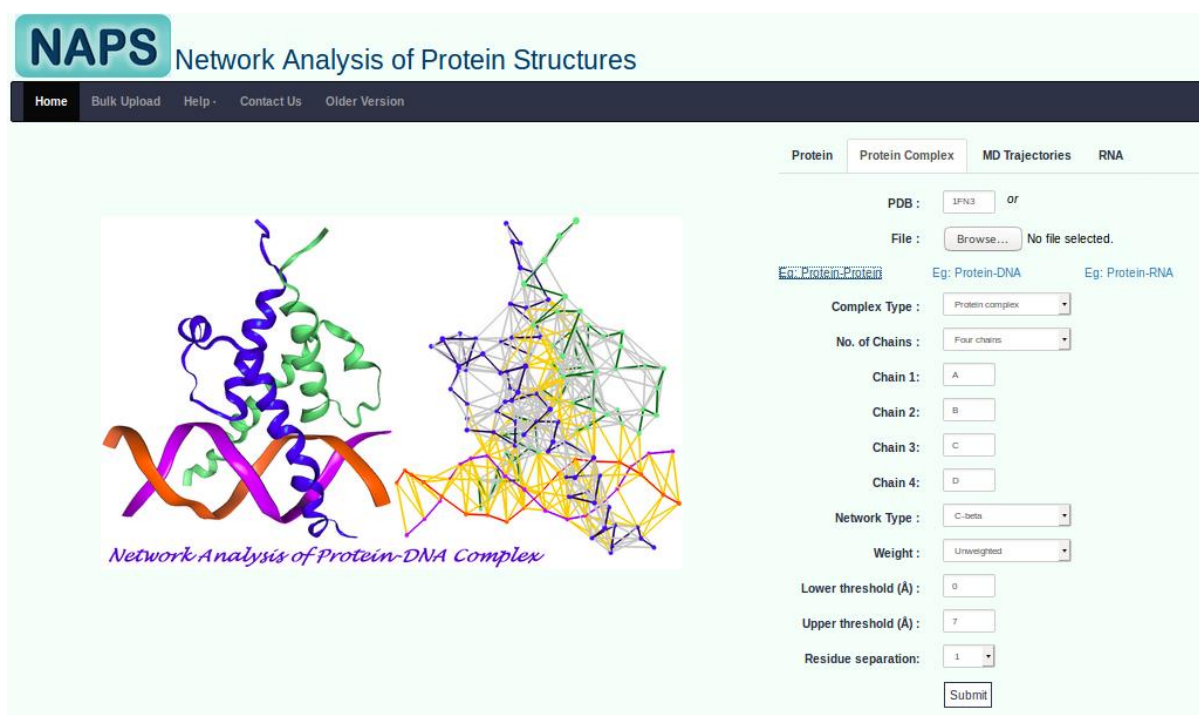


Figure 13.1: Input form for protein complex at the home page.

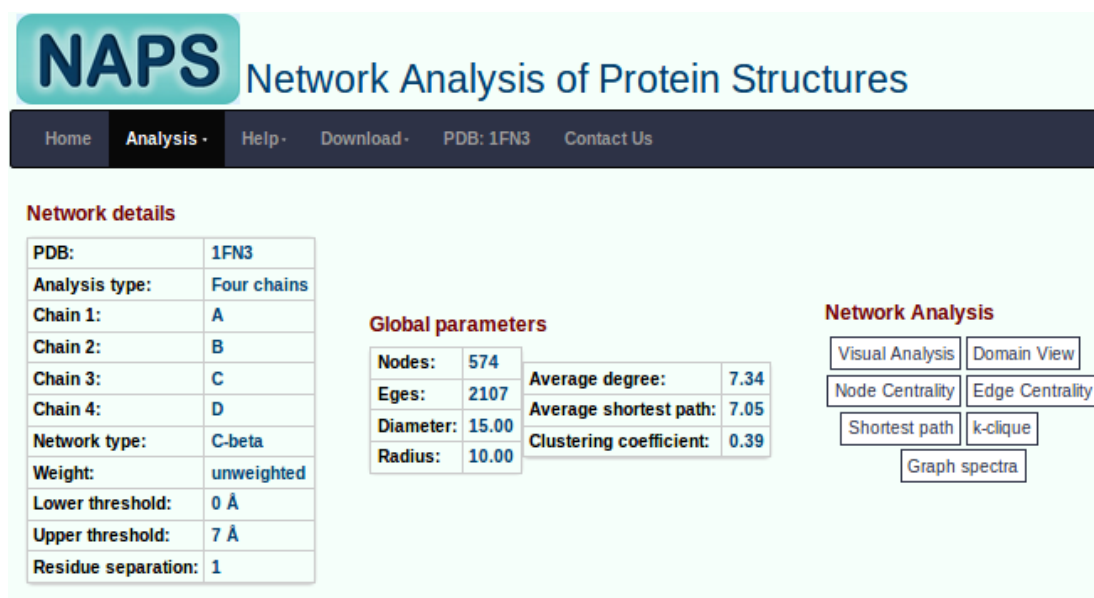


Figure 13.2: Network properties for protein complex.

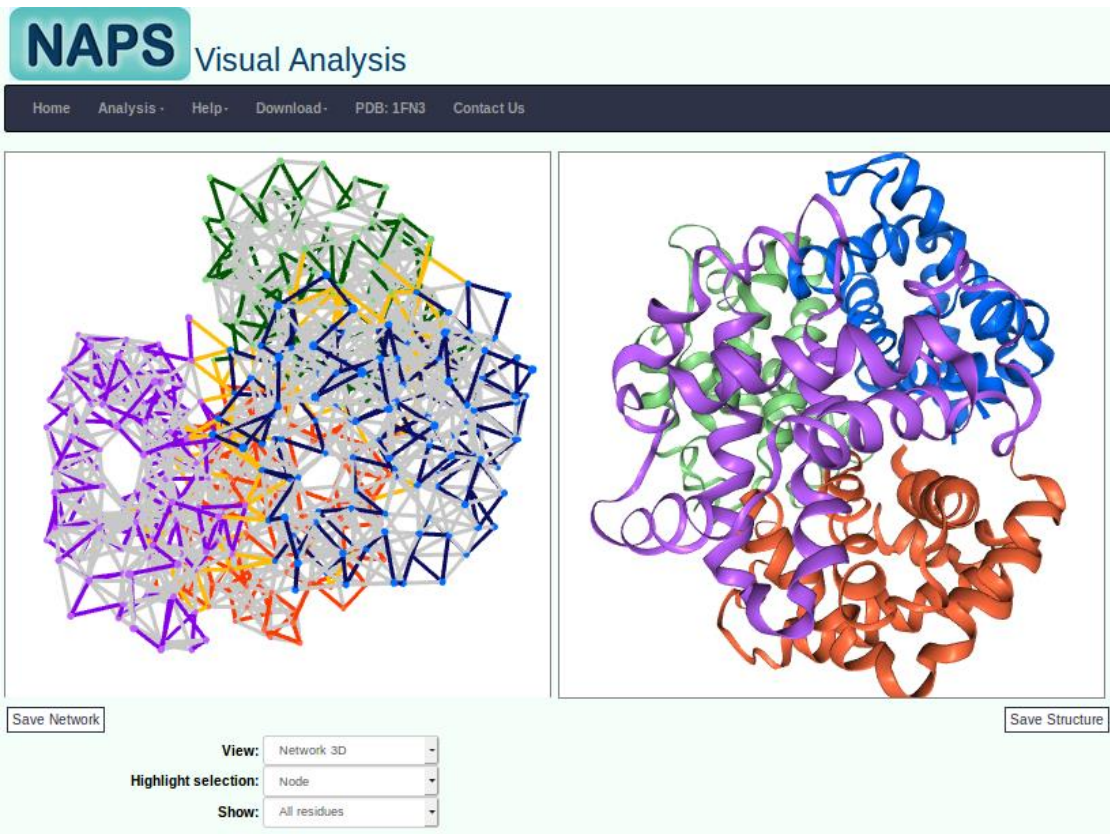


Figure 13.3: Network 3D view for protein complex. The two chains are shown in blue and green color in JSmol applet. The corresponding nodes in the network view are also shown in respective colors. Intra-chain edges are represented by grey color while inter-chain edges are represented by magenta color.

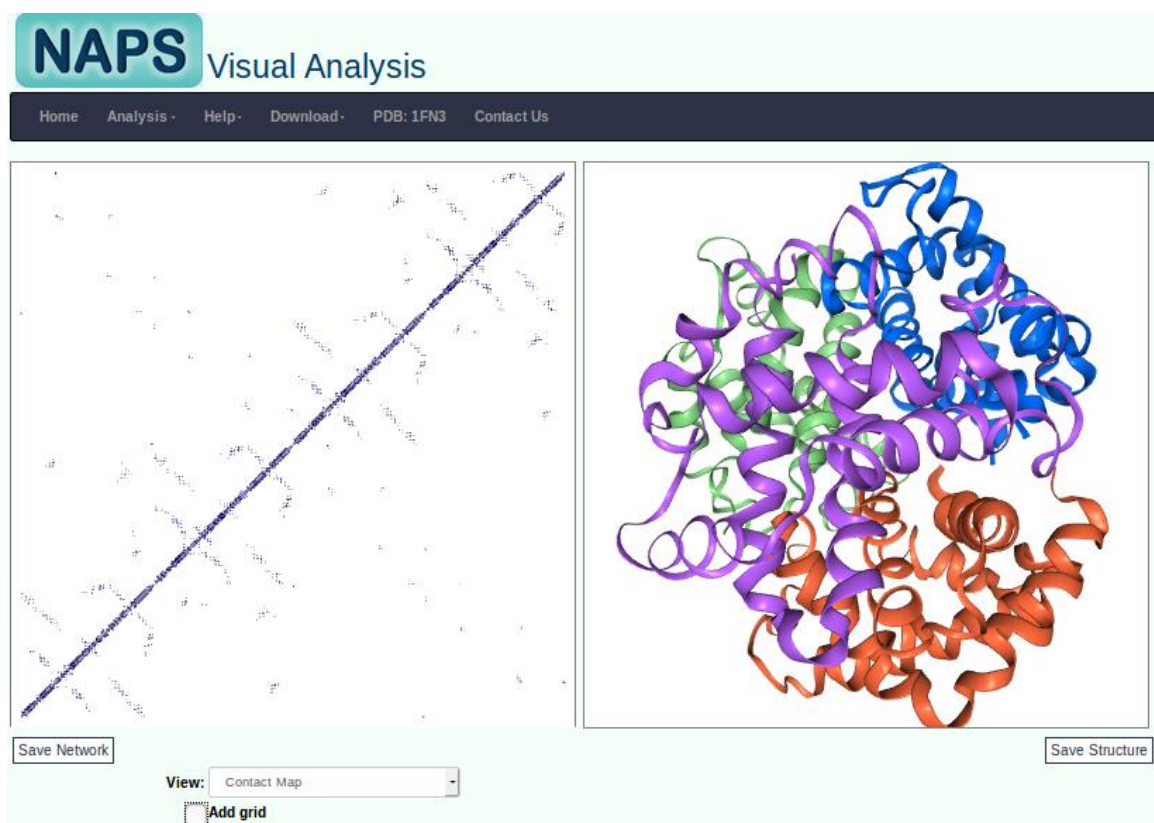


Figure 13.4: Contact map view of the protein complex.

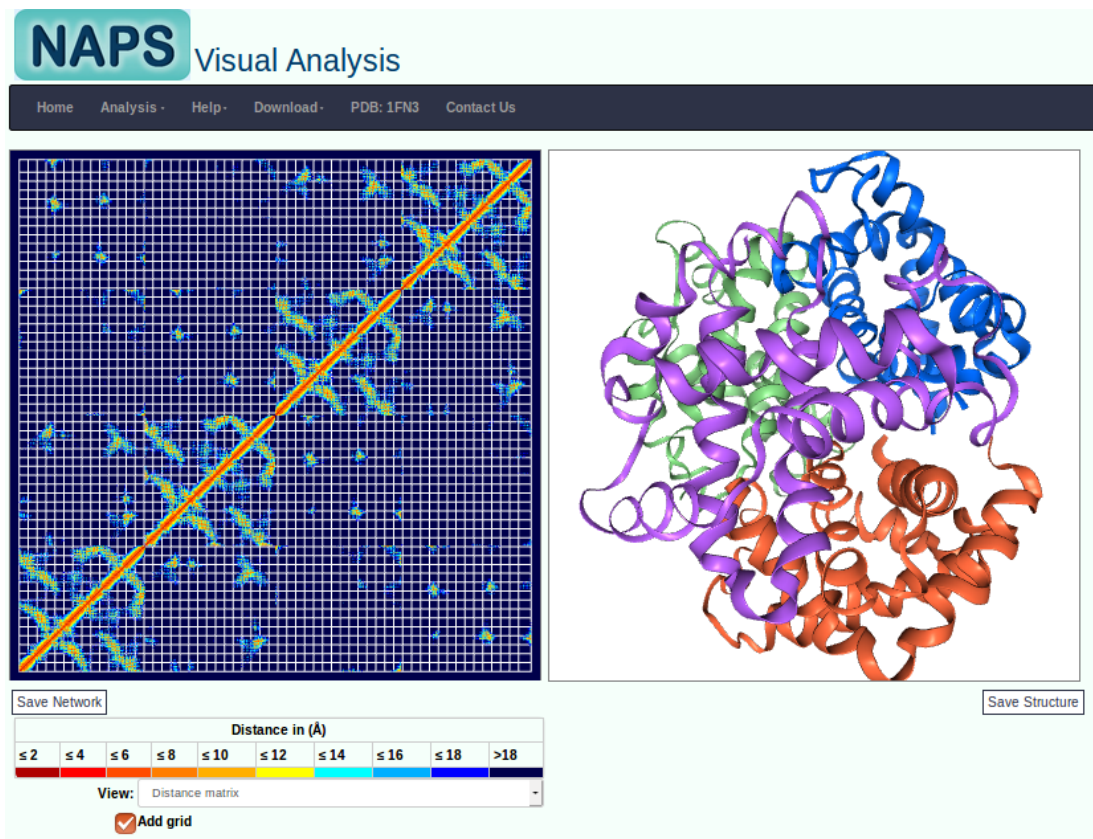


Figure 13.5: Distance matrix of the protein complex.

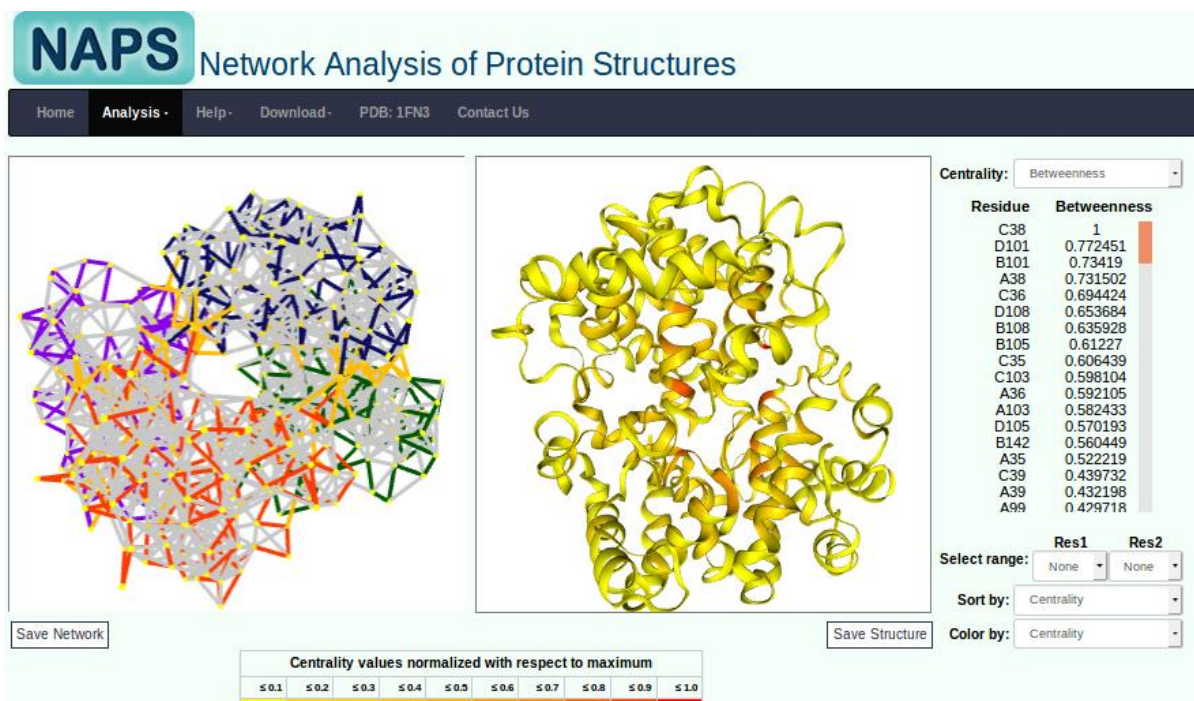


Figure 13.6: Betweenness centrality analysis of the protein complex. The interface residues with magenta edges showing protein-protein interaction can be observed to have high betweenness values indicating their importance in the inter-protein communications.

14. RNA Network

All the analysis that can be performed on one protein can be performed on a RNA structure. RNA structure network is constructed by considering a nucleotide residue as node and an edge is drawn between two nucleotide residues if any pairs of atoms of the two residues are within a threshold distance (default 5Å). This network representation is similar to atom pair contact network type of protein contact network. The 3D network view shows the nodes corresponding the nucleotide residues and backbone edges in blue color while all other non-backbone edges are shown in grey. The snapshots of the tool for different analysis pages are shown in following figures.

Figure 14.1: Input form for RNA structure at the home page.

Network details	
PDB:	2QBZ
Analysis type:	Single chain
Chain:	X
Network type:	Any atom
Weight:	unweighted
Lower threshold:	0 Å
Upper threshold:	5 Å
Residue separation:	1

Global parameters			
Nodes:	153	Average degree:	6.29
Edges:	481	Average shortest path:	5.86
Diameter:	18.00	Clustering coefficient:	0.47
Radius:	9.00		

Figure 14.2: Network properties for RNA.

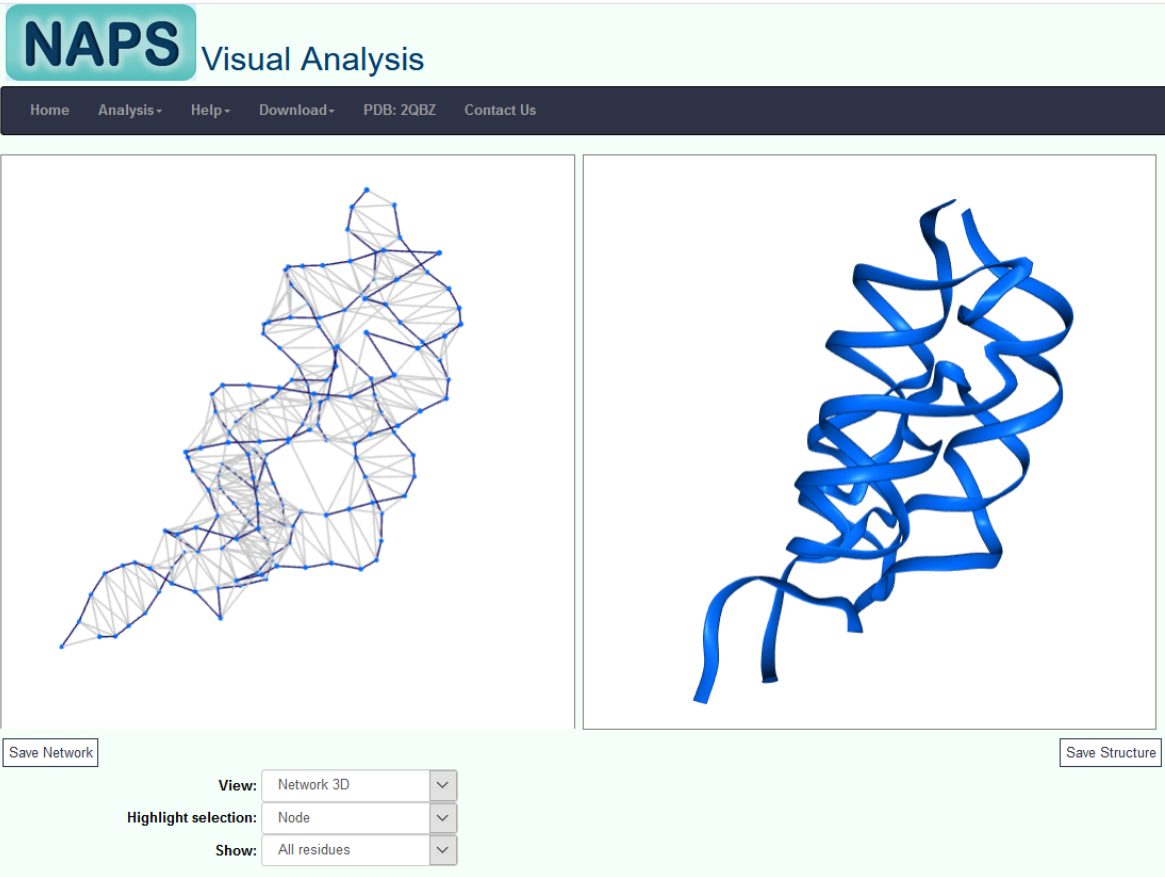


Figure 14.3: Network 3D view and molecular view for RNA.

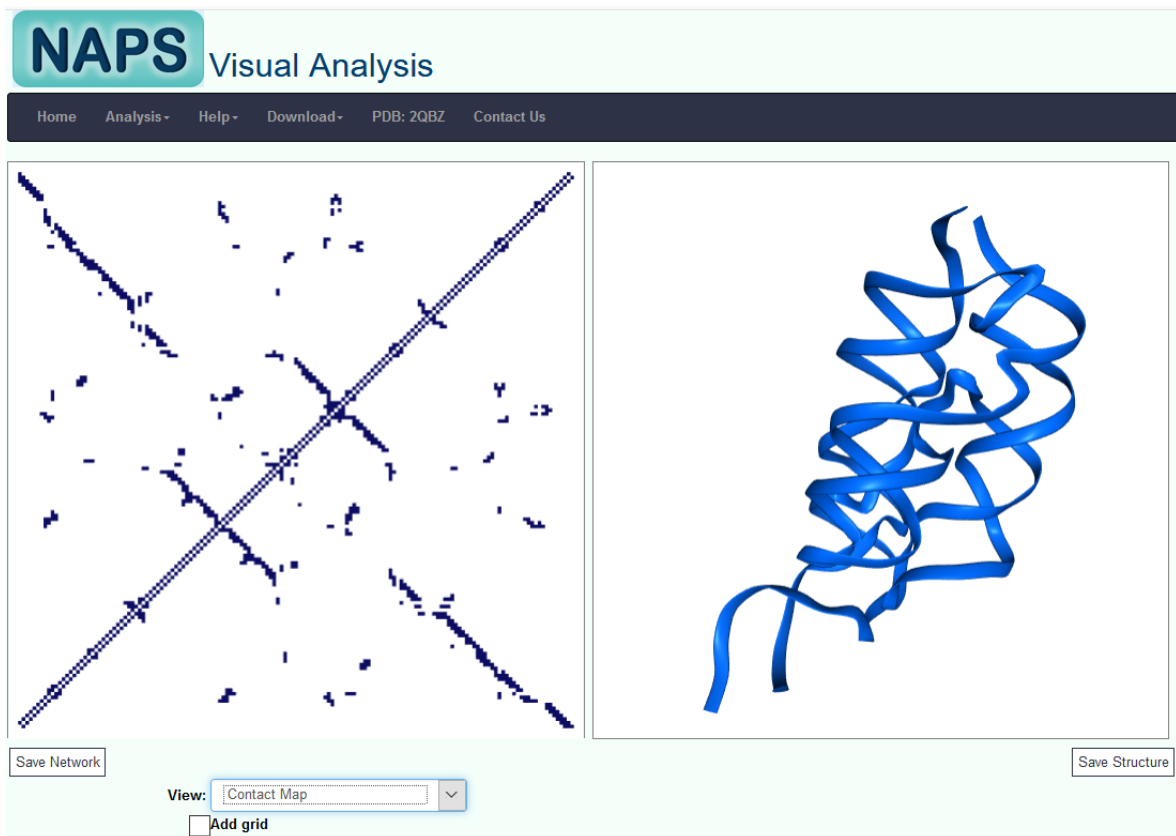


Figure 14.4: Contact map view of the RNA network.

NAPS Network Analysis of Protein Structures

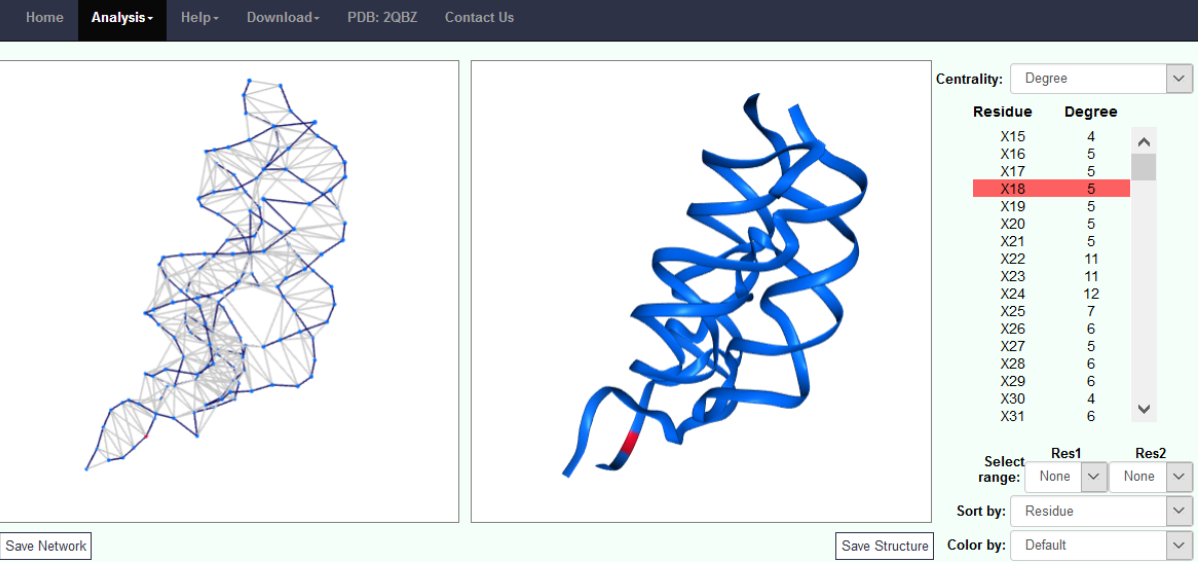


Figure 14.5: Degree centrality analysis of RNA network.

NAPS Shortest Path Analysis

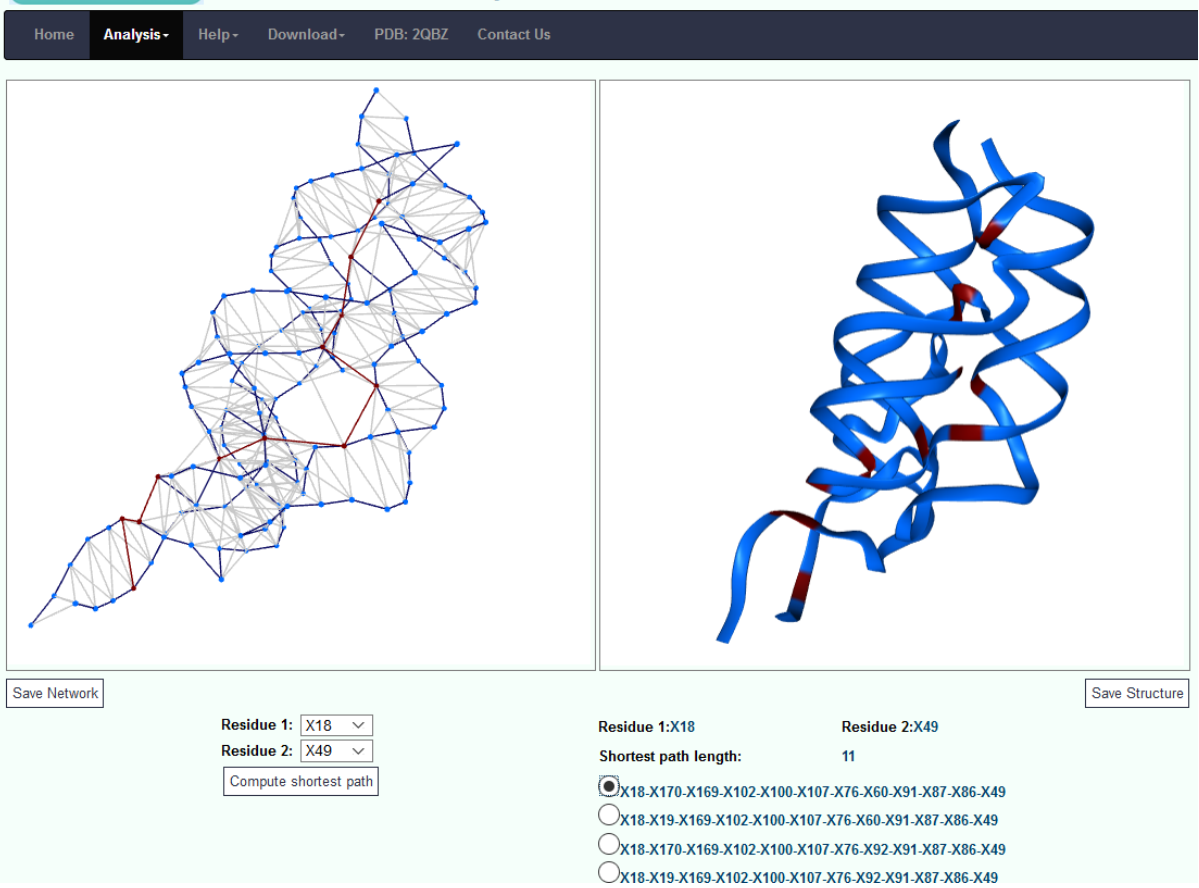


Figure 14.6: Shortest path analysis of RNA network.

15. Protein-Nucleic Acid Complexes

NAPS provides an option to analyze a protein structure in complex with a DNA or RNA. The takes the protein-nucleic acid complex structure in PDB format along with the network construction parameters as input. The input form is shown in Figure 15.1.

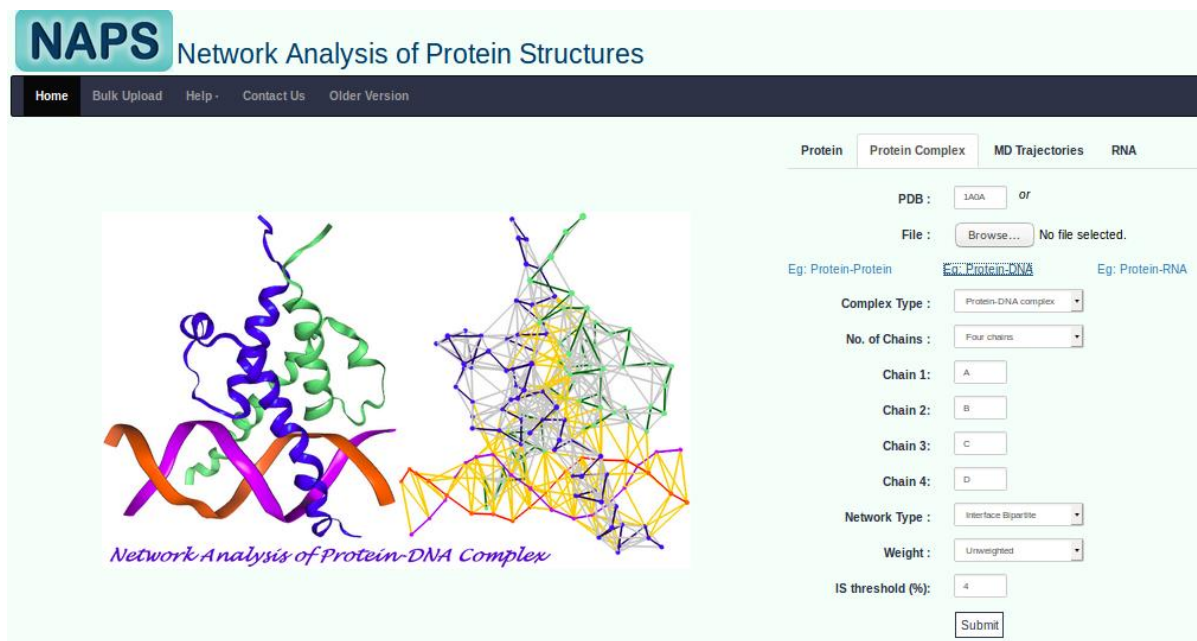


Figure 15.1: Input form for protein-nucleic acid complex.

A. Network Construction

- a) **Atom pair contact network:** Amino acid residues and nucleic acid residues in the complex are considered as nodes in the network and edges are constructed if the distance between any pair of atoms of the residue pair is within the lower and upper thresholds defined by the user (default upper threshold = 5 Å; lower threshold = 0 Å).

Unweighted: All edges are considered equally important.

Weighted: Edge weight is given by the number of atom pairs within cutoff distance.

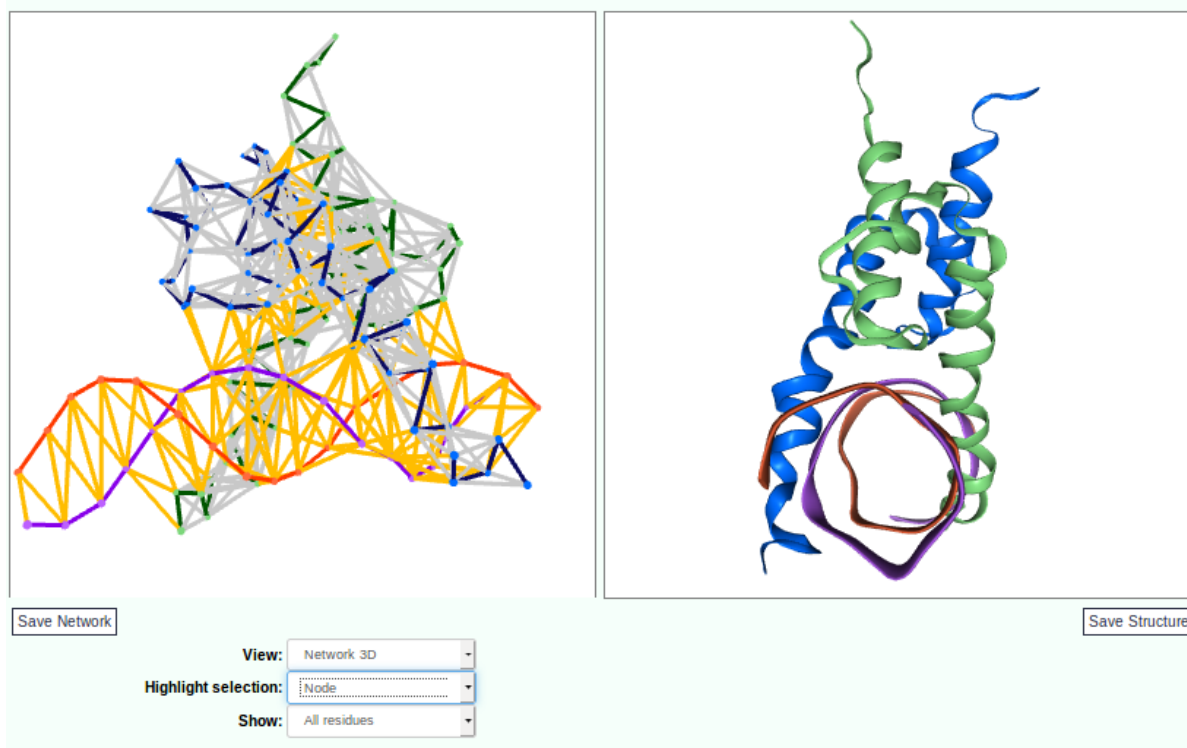


Figure 15.2: Network view and 3D molecular view of protein-DNA complex.

- b) **Interface bipartite network:** An amino acid residue and a nucleotide are considered as the two types of nodes in the bipartite network and an edge is constructed if the interaction strength between the two nodes is more than the threshold defined by the user (default = 4%). In this method, three types of interaction networks are constructed for the interaction of amino acids with the phosphate, sugar and base components of nucleotides respectively. The interaction strength as proposed by Sathyapriya, Vijayabaskar and Vishveshwara is calculated as:

$$I_{ij} = \left(\frac{n_{ij}}{\sqrt{N_i * N_j}} \right) * 100$$

where, n_{ij} is the number of side chain atom pairs of the nodes i and j within 4.5 Å. N_i and N_j are the normalization values of the residues i and j given by Kannan and Vishveshwara as shown in Table 15.1.

Unweighted: All edges are considered equally important.

Weighted: The interaction strength (I_{ij}) is considered as edge weight. The choice of the threshold depends on the biological problem to be addressed.

Table 15.1: Normalization value used to construct Interaction Strength Network of protein-nucleic acids.

Residue Type	Normalization value
Ala	55.76
Arg	93.79
Asn	73.41
Asp	75.15
Cys	54.95
Gln	78.13
Glu	78.83
Gly	47.31
His	83.74
Ile	67.95
Leu	72.25
Lys	69.61
Met	69.26
Phe	93.31
Pro	51.33
Ser	61.39
Thr	63.71
Trp	106.70
Tyr	100.72
Val	62.37

Protein-DNA complex

Nucleotide Component	Normalization value
Phosphate	25
Sugar	28
Adenine	144
Guanine	156
Cytosine	114
Thymine	120

Protein-RNA complex

Nucleotide Component	Normalization value
Phosphate	28
Sugar	49
Adenine	91
Guanine	97
Cytosine	59
Uracil	66

Connected component analysis: A connected component (or cluster) of a graph is a subgraph in which any two vertices are connected to each other by paths, i.e. any all nodes can be reached from all other nodes of the subgraph. The connected components of protein-phosphate bipartite network are show in Figure 15.3.

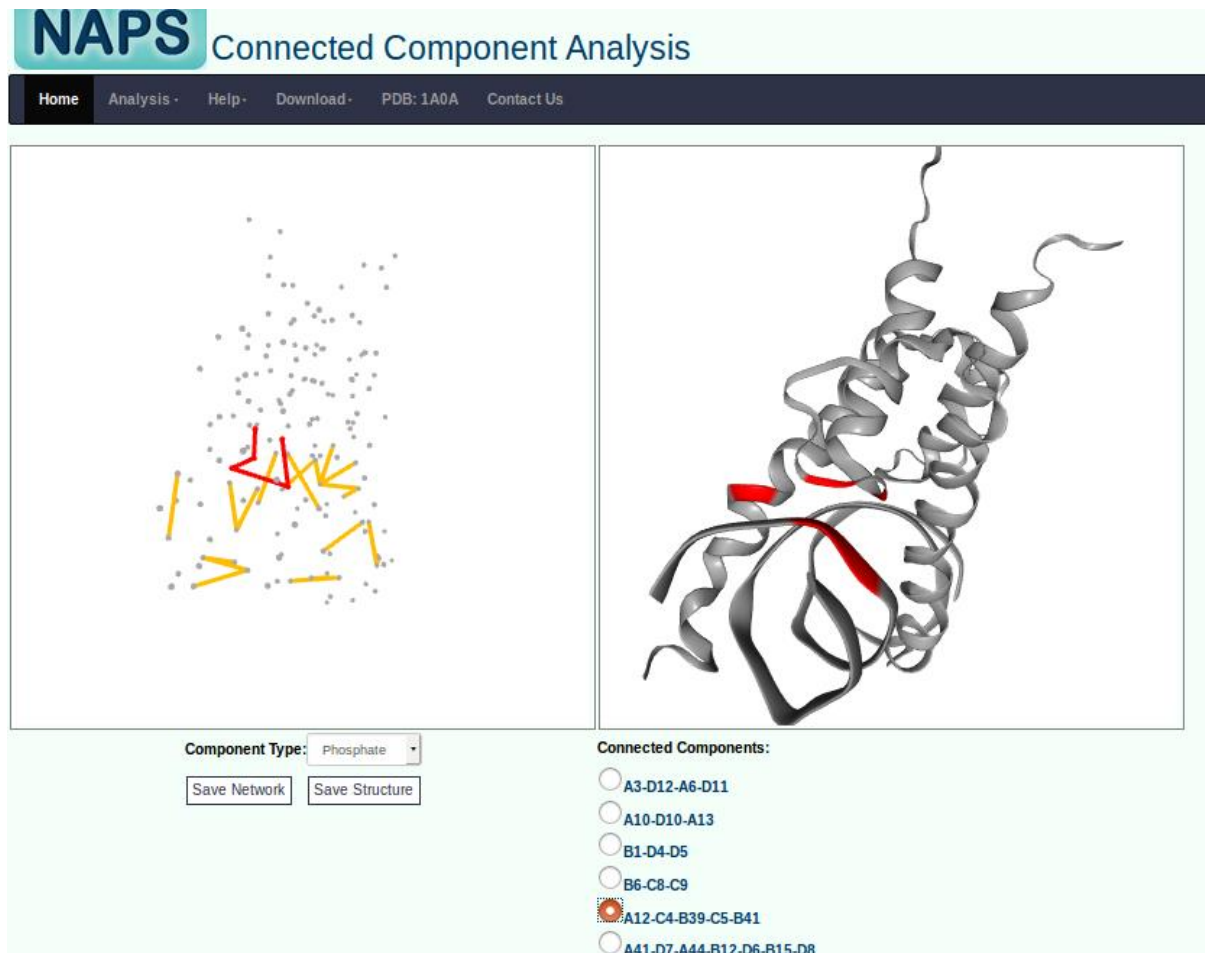


Figure 15.3: Connected components analysis.

Hub analysis: An amino acid residue is considered as a hub if the degree of the node is \geq threshold (default 4). The hubs in a protein-base bipartite network are highlighted by red color in Figure 15.4.

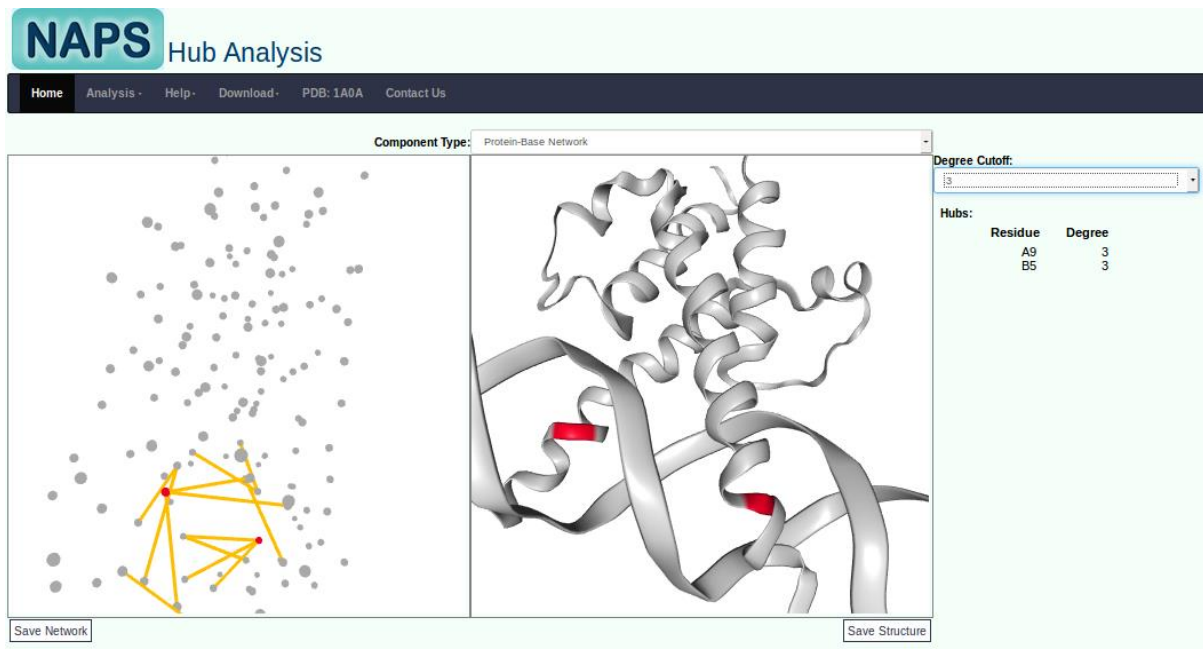


Figure 15.4: Hubs are highlighted in a protein-base bipartite network.

16. Network Analysis of Molecular Dynamics Data

The portal provides an option to analyze Molecular Dynamics (MD) trajectory data, which can be uploaded in DCD format. The basic parameters related to trajectory such as starting timestep, ending timestep and the stride value (time interval) can be selected along with the network type and its parameters. The stride value is used to stride over frames included in the simulation trajectory. For example, if there are 1000 frames in the trajectory, specifying a stride value of 10 will consider every 10th frame in the trajectory.

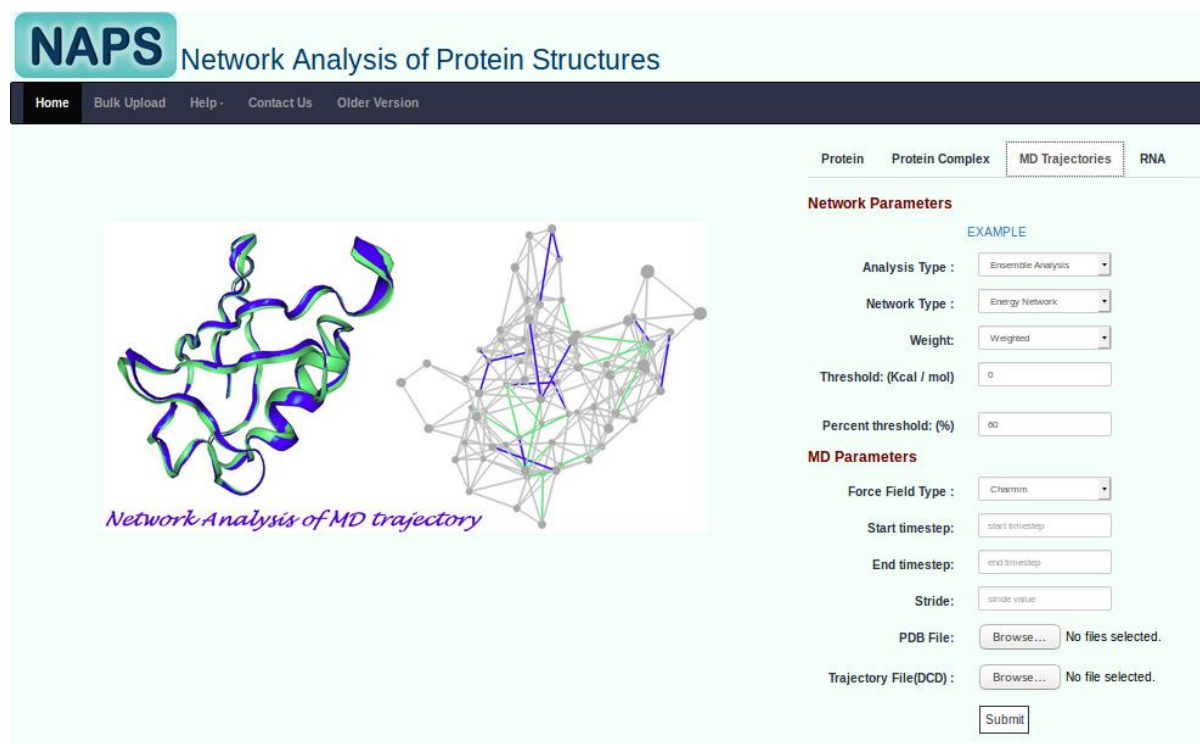


Figure 16.1: Input form for network analysis of molecular dynamics trajectory at the home page.

A. Dynamic Cross Correlation Map

The dynamic cross correlation is a measure of similarity of two amino acid residues as a function of the displacement of one relative to the other. The cross-correlation between the displacements of the residues along the trajectory is calculated as follows:

$$DCC = \frac{\langle \Delta r_i(t) \cdot \Delta r_j(t) \rangle_t}{\sqrt{\langle \|\Delta r_i(t)\|^2 \rangle_t} \sqrt{\langle \|\Delta r_j(t)\|^2 \rangle_t}}$$

where $r_i(t)$ denotes the vector of the i^{th} atom's coordinates as a function of time t . $\langle . \rangle_t$ means the time ensemble average and $\Delta r_i(t) = r_i(t) - \langle r_i(t) \rangle_t$.

The users can view the cross correlation map over the entire trajectory between all the residue pairs. This is present along with options to download the correlation map and matrix. The DCC for an example trajectory is shown in Figure 16.2.

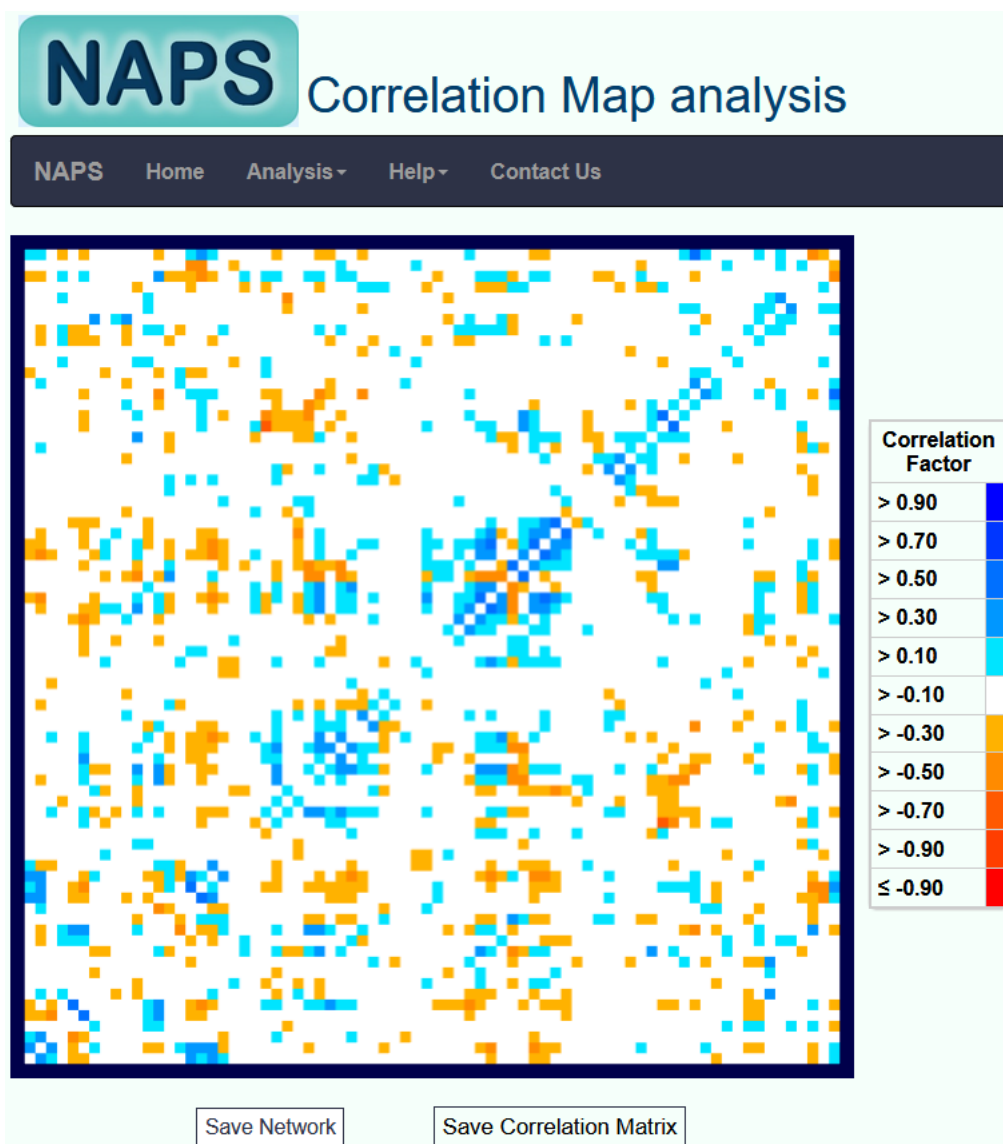


Figure 16.2: Dynamic Cross Correlation map

B. Network Construction

The portal provides options to construct the following four types of networks for MD data:

- C-alpha network:** For each timestep, an amino acid residue represented by the C-alpha atom is considered as node in the network and an edge is constructed if the distance between a pair of C-alpha atoms is within the lower and upper thresholds defined by the user (default upper threshold = 7 Å; lower threshold = 0 Å).

Unweighted: All edges are considered equally important.

Weighted: Edge weight for a C-alpha weighted network is given by:

$$w_{ij} = \frac{1}{d_{ij}}$$

where d_{ij} is the euclidean distance between C-alpha atoms of i^{th} and j^{th} residues.

- b) **C-beta network:** For each timestep, an amino acid residue represented by the C-beta atom is considered as node in the network and an edge is constructed if the distance between the C-beta atoms (C-alpha for GLY) is within the lower and upper thresholds defined by the user (default upper threshold = 7 Å; lower threshold = 0 Å).

Unweighted: All edges are considered equally important.

Weighted: Edge weight weighted network is given by:

$$w_{ij} = \frac{1}{d_{ij}}$$

where d_{ij} is the euclidean distance between C-beta atoms of i^{th} and j^{th} residues.

- c) **Energy network:** For each timestep, an amino acid represented by the C-alpha atom is considered as node and an edge is constructed if the interaction energy between residues lies between the given cutoff (default upper threshold = 1 Kcal/mol, default lower threshold = 0.1 Kcal/mol). The interaction energy between any 2 given residues primarily consists of 2 contributions - The Van der Waal interaction energy and the electrostatic interaction energy, which is give by:

$$Energy = \frac{\sum_{ij}(Electrostatic\ potential\ \langle i, j \rangle + Lennard\ Jones\ Potential\ \langle i, j \rangle)}{N\langle i, j \rangle}$$

where i and j are the individual atoms of the two residues and $N\langle i, j \rangle$ = no of $\langle i, j \rangle$ pairs between the two residues.

Unweighted: All edges are considered equally important.

Weighted: Edge weight of network is equal to the normalized interaction energy:

$$w_{ij} = (-1 * Energy) + C$$

where C is a large positive normalization constant

- d) **Cross Correlation Network:** A network is constructed for the entire trajectory based on the cross correlation of residues. An amino acid residue represented by the C-alpha atom is considered the node and an edge is constructed if the absolute value of the cross correlation between residues is within the user specified cutoffs (lower cutoff default = 0.5 and upper cutoff default = 1).

Unweighted: All edges are considered equally important.

Weighted: Edge weight of network is equal to the measure of the absolute cross correlation value over the trajectory between the residues

C. Centrality Analysis

Centrality measure of a node provides a quantification of the topological importance of the node in the network. All the centrality analysis options provided in NAPS for a single protein structure are available for MD trajectory analysis. The following centrality analysis options for MD analysis are provided in the NAPS portal:

- a) **Centrality Analysis of a particular timestep of the trajectory:** One timestep of the trajectory can be selected from the dropdown list and extensive centrality analysis can be carried out for it. A snapshot of the centrality analysis for a timestep is shown in Figure 16.3.

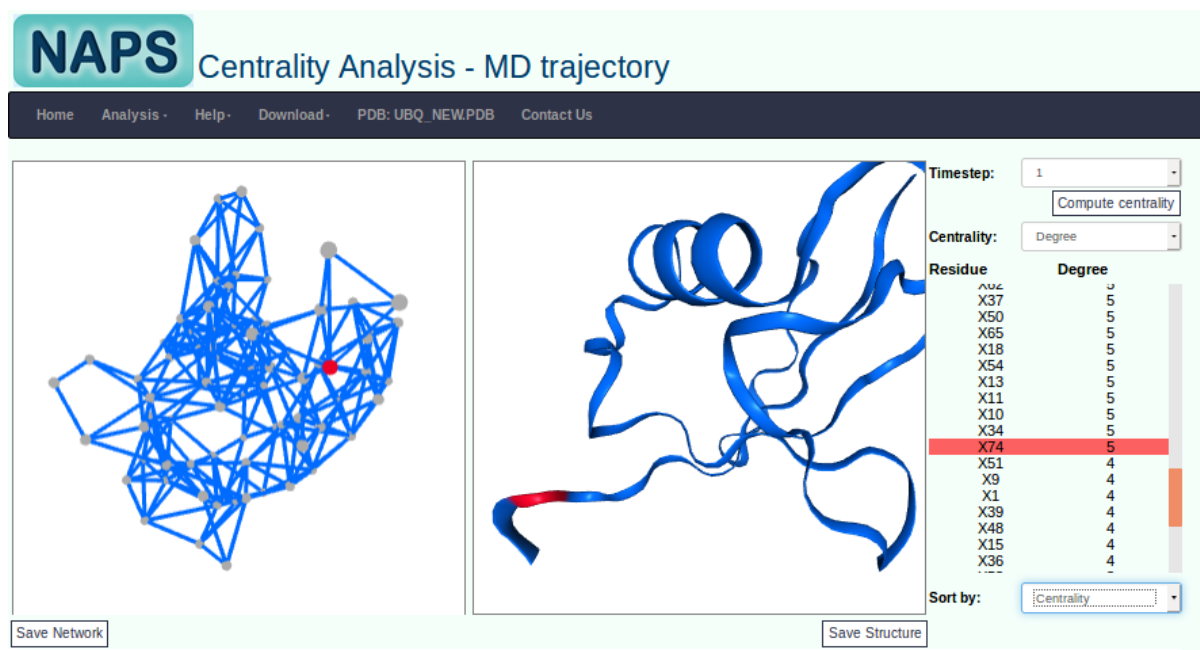


Figure 16.3: Centrality analysis for a timestep.

- b) **Comparative centrality analysis of any two timesteps of the trajectory:** Two timesteps of the trajectory can be selected from the dropdown lists along with centrality measure to be analyzed. The selected centrality measure can be compared for all the residues of the protein (or complex). An interactive plot of the centrality measure is also provided to for ease of analysis. A snapshot of the comparative analysis is shown in Figure 16.4.

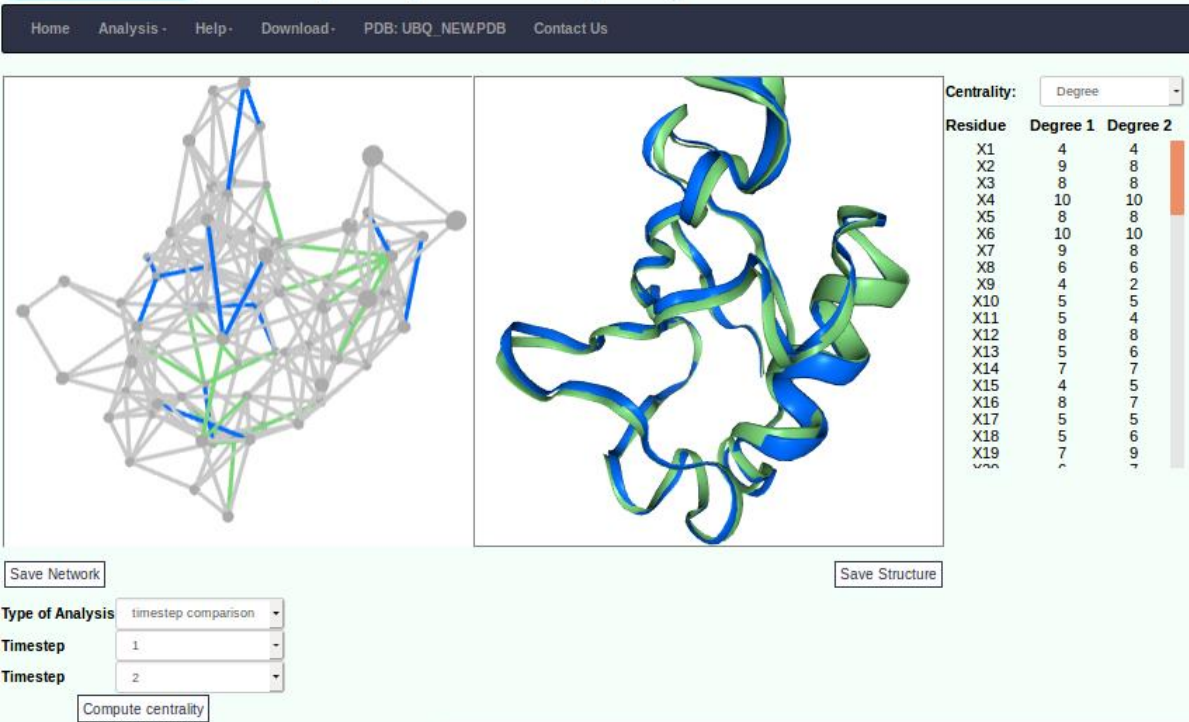


Figure 16.4: Comparative centrality analysis of two timesteps.

- c) **Comparative centrality analysis to the average centrality measures over the trajectory:** The centrality of a particular timestep can be compared with the average centrality of all the timesteps across the trajectory. An interactive plot of the centrality measure is also provided to for ease of analysis. A snapshot of the comparative analysis is shown in Figure 16.5.

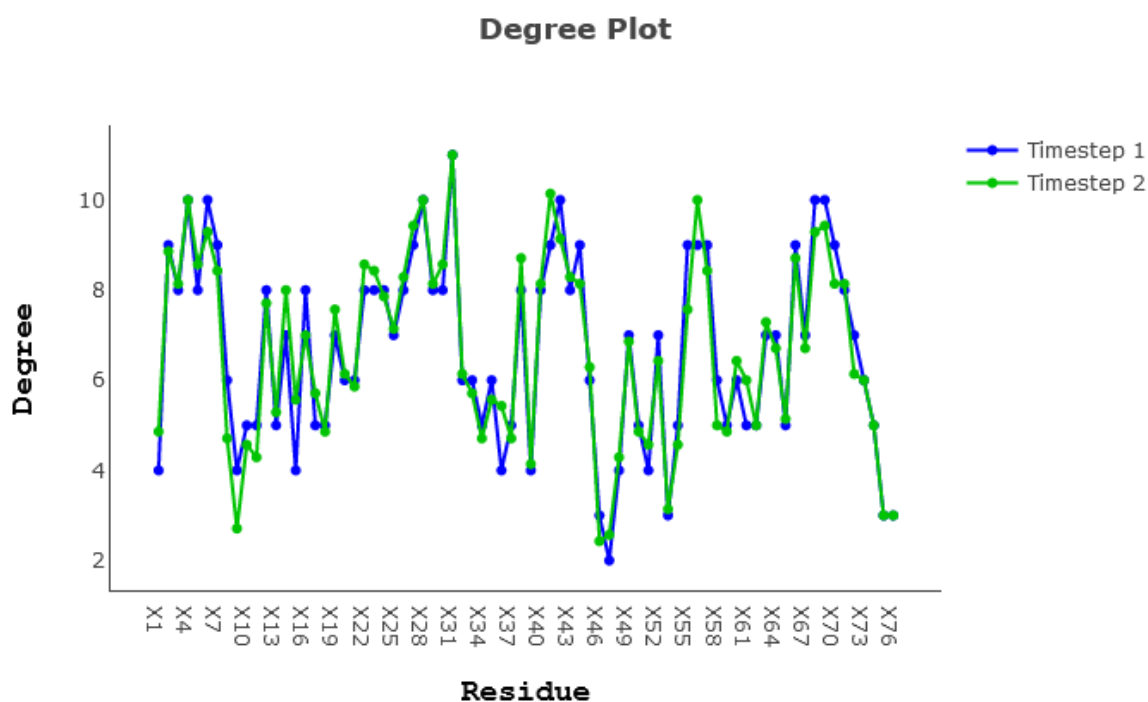


Figure 16.5: Plot of Degree centrality of a timestep overlapped with the average Degree across all the timesteps of the trajectory.

D. Shortest Path Analysis

The shortest path analysis discussed earlier for a static protein structure can be carried out for the MD data also. The following shortest path analysis are available in NAPS portal:

- a) **Shortest path analysis of a particular timestep:** One timestep of the trajectory can be selected from the dropdown list and all shortest paths can be obtained for a particular pair of selected residues. The snapshot of shortest path analysis for a timestep is shown in Figure 16.6.

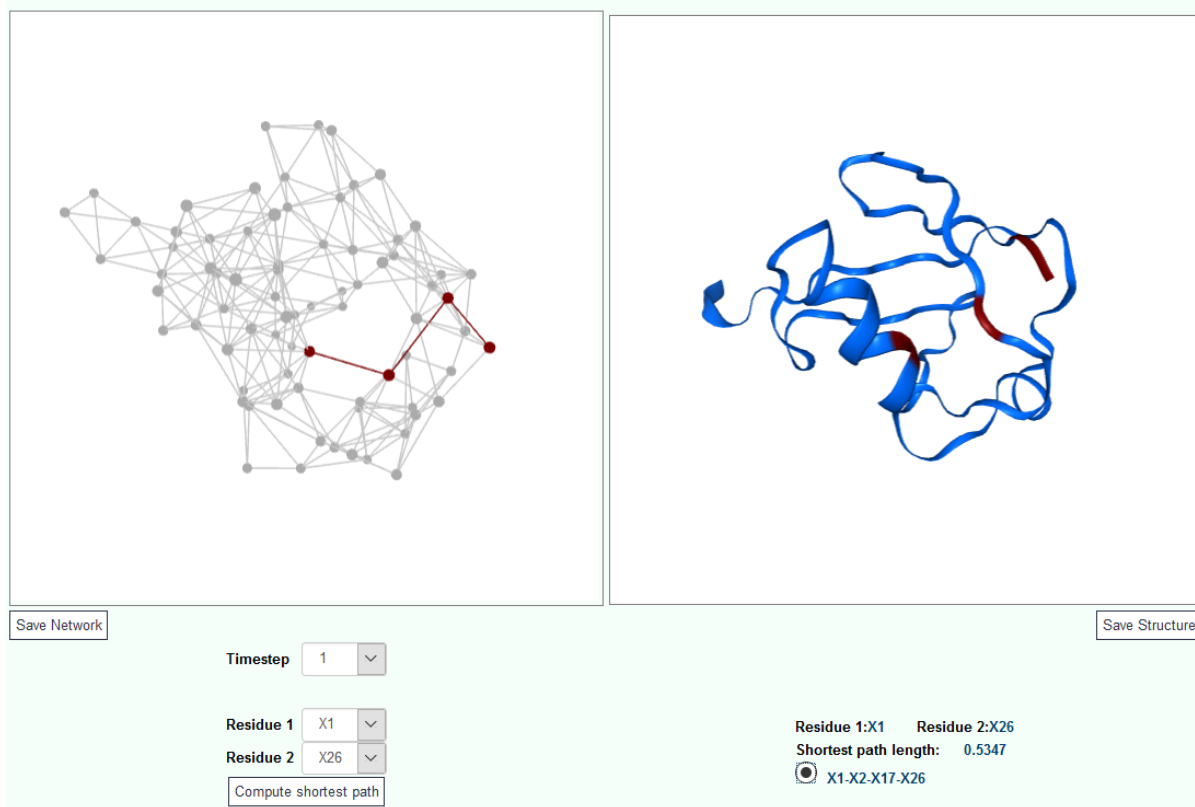


Figure 16.6: Shortest path analysis of a timestep.

- b) **Comparative shortest path analysis of any two timesteps of the trajectory:** Two timesteps of the trajectory can be selected from the dropdown lists along with pair of residues for which shortest path analysis has been carried out. The portal provides all the shortest paths for the two timesteps. The list of common paths are shown differently from the list of unique paths in each of the timesteps.

17. References

1. Kannan,N. and Vishveshwara,S. (1999) Identification of side-chain clusters in protein structures by a graph spectral method. *J. Mol. Biol.*, **292**, 441–464.
2. Brinda,K.V. and Vishveshwara,S. (2005) A network representation of protein structures: implications for protein stability. *Biophys. J.*, **89**, 4159–4170.
3. Bagler,G. and Sinha,S. (2005) Network properties of protein structures. *Physica A: Statistical Mechanics and its Applications*, **346**, 27–33.
4. Bartoli,L., Fariselli,P. and Casadio,R. (2007) The effect of backbone on the small-world properties of protein contact maps. *Phys Biol*, **4**, L1–5.
5. Bagler,G. and Sinha,S. (2007) Assortative mixing in Protein Contact Networks and protein folding kinetics. *Bioinformatics*, **23**, 1760–1767.
6. Vishveshwara,S., Ghosh,A. and Hansia,P. (2009) Intra and inter-molecular communications through protein structure network. *Curr. Protein Pept. Sci.*, **10**, 146–160.
7. Patra,S.M. and Vishveshwara,S. (2000) Backbone cluster identification in proteins by a graph theoretical method. *Biophys. Chem.*, **84**, 13–25.
8. Chakrabarty,B. and Parekh,N. (2012) Analysis of graph centrality measures for identifying Ankyrin repeats. In *2012 World Congress on Information and Communication Technologies (WICT)*.pp. 156–161.
9. Chakrabarty,B. and Parekh,N. (2014) Identifying tandem Ankyrin repeats in protein structures. *BMC Bioinformatics*, **15**, 6599.
10. Chakrabarty,B. and Parekh,N. (2014) PRIGSA: protein repeat identification by graph spectral analysis. *J Bioinform Comput Biol*, **12**, 1442009.
11. Atilgan,A.R., Akan,P. and Baysal,C. (2004) Small-world communication of residues and significance for protein dynamics. *Biophys. J.*, **86**, 85–91.
12. Greene,L.H. and Higman,V.A. (2003) Uncovering network systems within protein structures. *J. Mol. Biol.*, **334**, 781–791.
13. Aftabuddin,M. and Kundu,S. (2007) Hydrophobic, Hydrophilic, and Charged Amino Acid Networks within Protein. *Biophys J*, **93**, 225–231.
14. Alves,N.A. and Martinez,A.S. (2007) Inferring topological features of proteins from amino acid residue networks. *Physica A: Statistical Mechanics and its Applications*, **375**, 336–344.
15. JSmol: an open-source HTML5 viewer for chemical structures in 3D.
16. Yan,W., Zhou,J., Sun,M., Chen,J., Hu,G. and Shen,B. (2014) The construction of an amino acid network for understanding protein structure and function. *Amino Acids*, **46**, 1419–1439.
17. Kyte,J. and Doolittle,R.F. (1982) A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.*, **157**, 105–132.

