



ISSN: 0067-2904

## Heuristic Modularity for Complex Identification in Protein-Protein Interaction Networks

Amenah H. Abdulateef<sup>1,3</sup>, Bara'a A. Attea\*<sup>2</sup>, Ahmed N. Rashid<sup>3</sup>

<sup>1</sup>Department of Computer Science, College of Education for Pure Science (Ibn al-Haitham), University of Baghdad, Iraq

<sup>2</sup>Department of Computer Science, College of Science, University of Baghdad, Baghdad, Iraq

<sup>3</sup>Computers and Information Technology, University of Anbar, Anbar, Iraq

### Abstract

Due to the significant role in understanding cellular processes, the decomposition of Protein-Protein Interaction (PPI) networks into essential building blocks, or complexes, has received much attention for functional bioinformatics research in recent years. One of the well-known bi-clustering descriptors for identifying communities and complexes in complex networks, such as PPI networks, is modularity function. The contribution of this paper is to introduce heuristic optimization models that can collaborate with the modularity function to improve its detection ability. The definitions of the formulated heuristics are based on nodes and different levels of their neighbor properties. The modularity function and the formulated heuristics are then injected into the mechanism of a single objective Evolutionary Algorithm (EA) tailored specifically to tackle the problem, and thus, to identify possible complexes from PPI networks. In the experiments, different overlapping scores are used to evaluate the detection accuracy in both complex and protein levels. According to the evaluation metrics, the results reveal that the introduced heuristics have the ability to harness the accuracy of the existing modularity while identifying protein complexes in the tested PPI networks.

**Keywords:** complex detection; graph partitioning; heuristic; modularity; Protein-Protein Interaction network.

### دالة النمطية الموجهة لتحديد المركبات في شبكات التفاعل بين البروتين والبروتين

امنة هيثم<sup>1,3</sup>، براء علي عطية\*<sup>2</sup>، احمد رشيد<sup>3</sup>

<sup>1</sup>قسم علوم الحاسوب، كلية التربية للعلوم الصرفة (ابن الهيثم)، جامعة بغداد، العراق

<sup>2</sup>قسم علوم الحاسوب، كلية العلوم، جامعة بغداد، بغداد، العراق

<sup>3</sup>أجهزة الكمبيوتر وتكنولوجيا المعلومات، جامعة الأنبار، الأنبار، العراق

### الخلاصة

نظرًا للدور المهم في فهم العمليات الخلوية، فإن تحليل شبكات التفاعل بين البروتين والبروتين (PPI) إلى لبنات أساسية أو مجتمعات (مركبات)، قد حظي في السنوات الأخيرة إلى الكثير من الاهتمام في بحوث المعلوماتية الحيوية. إحدى الدوال الرياضية المعروفة للكشف عن المجتمعات في الشبكات المعقدة، مثل شبكات ال PPI، هي دالة النمطية. تتمثل مساهمة هذا البحث في تقديم نماذج تحسين موجهة استكشافية التي يمكن أن تتعاون مع دالة النمطية التقليدية لتحسين قدرتها على تحديد المركبات. تعتمد تعريفات نماذج التحسين الموجهة المقترحة في هذا البحث على العقد (أي البروتينات) ومستويات مختلفة من خصائص

الجيران لهذه البروتينات. يتم بعد ذلك حقن دالة النمطية و نماذج التحسين الموجهة المصاغة في آلية خوارزمية تطويرية (EA) مصممة خصيصًا لمعالجة المشكلة ، وبالتالي ، لتحديد المركبات المحتملة من شبكات ال PPI. في التجارب، يتم استخدام درجات تداخل مختلفة لتقييم دقة الكشف في كل من المستويات المعقدة والبروتين. وفقًا لمقاييس التقييم، تكشف النتائج أن نماذج التحسين الموجهة التي تم إدخالهم إلى خوارزمية ال (EA) لديهم القدرة على تسخير دالة النمطية التقليدية مع تحديد مركبات البروتين في شبكات ال PPI المختبرة.

## 1. Introduction

Many complex systems in all areas of science, including social science, politics, biology and medicine, can be represented as networks. Topological analyses of such complex networks are universal and provide insights in many science studies. Complex systems are usually organized in compartments, which have their own role and / or function. In the network representation, such compartments appear as sets of nodes with a high density of internal links, whereas links between compartments have a lower density. These subgraphs are called communities, or modules, and can occur in a wide variety of networked systems. Finding compartments may shed light on the organization of complex systems and on their function. Therefore, detecting communities in networks has become a fundamental problem in network science. Many methods have been developed, using tools and techniques from different disciplines like physics, applied mathematics, biology, computer and social sciences. However, it is still not clear which algorithms are reliable and shall to be used in applications [1].

As an increasing amount of protein–protein interaction (PPI) data becomes available, its computational interpretation has become an important problem in bioinformatics. Observations show that PPI networks possess invaluable evolutionary insights and information to understand various biological processes and cellular functions. However, prediction of protein complexes, like many other practical optimization problems, falls into the category of strongly NP-hard combinatorial optimization problems that can easily bewilder exact optimization algorithms [2], [3].

Complex network clustering is data clustering in dividing the interested entities into clusters or modules. However, clusters in complex networks are based on both the inter and intra connections densities, while clusters in data clustering are groups of points close to each other in a way forming several local optima. In the literature, modularity [3] and a series of follow-up models have been proposed to measure the quality of a set of predicted intra-dense and inter-sparse subgraphs in a graph. The majority of these works have been applied to community detection in social networks and to complex identification in PPI networks.

The main contribution of this paper is to develop a heuristic ground definitions for modularity that can improve its detection ability. Here we are motivated by the logical consequence of protein neighboring properties and how to exploit and couple such properties with modularity function. In this paper, a single objective Evolutionary Algorithm (EA) [4] is adopted to combine modularity (as an objective function) and the proposed heuristic approaches. The algorithm attempts, with aid of modularity, to identify the *global structure* of the complexes and with the aid of heuristic functions, to fine tune such complexes. In the experimental results, we show that coupling the proposed heuristic operator as exploiter to capture local structures of the solutions provided by modularity can significantly improve the detection performance of EA.

In the remainder of this paper, preliminary concepts relating to complex detection problem in PPI networks and the main interest in the literature towards solving complex detection problem in PPI networks are presented first. These are followed by a closer look into the formal development of the proposed evolutionary based complex detection framework together with the proposed heuristic operators. Experimental results are then provided to support the positive impact of the proposed heuristic definitions to further correct the complex structures of the well-known modularity function. The final section of this paper presents major conclusions and further directions of this work.

## 2. Related background

A complex network such as PPI network can be denoted by  $\mathcal{N}$ , of  $n$  proteins and  $m$  interactions. In other terms  $\mathcal{N}$  is said to be with cardinality  $n$  and volume  $m$ . Generally,  $\mathcal{N}$  can be modeled as undirected graph  $G = (V, E)$  of a set  $V = \{v_1, v_2, \dots, v_n\}$  of  $n$  vertices, and a set  $E \subseteq V \times V = \{(v_i, v_j) | v_i, v_j \in V \text{ and } i \neq j\}$  of  $m$  edges. Note that throughout this paper, the terms: tie, edge, link,

connection, relation, and interaction are used interchangeably to denote any vertex pair  $(v_i, v_j)$  in  $E$ . Also, let  $\Omega$  be the space of all possible partitioning solutions for  $\mathcal{N}$  and let  $\mathcal{C} \in \Omega = \{C_1, \dots, C_K\}$  be a network partitioning solution belongs to the space  $\Omega$  with  $K$  partitions or divisions. Normally, any unsigned graph  $G$  can be represented by a symmetric  $n \times n$  adjacency matrix denoted by  $A$ . Rows and columns of  $A$  are labeled with the vertices of  $V$  and assigned with 1 in entry  $(i, j)$  if vertex pair  $(v_i, v_j)$  is in  $E$ , and set to 0 if  $(v_i, v_j) \notin E$ .

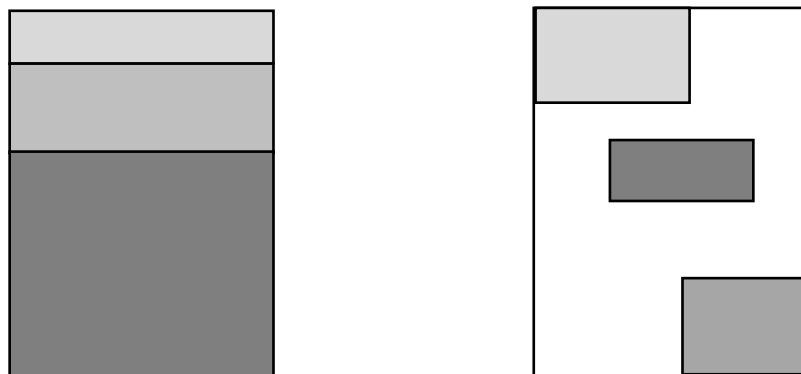
From the adjacency matrix  $A$  a set of  $n$  direct neighboring lists,  $L = \{l_1, l_2, \dots, l_n\}$ , can be formed. Each list  $l_i$  in the set  $L$  aggregates connections of all vertices with vertex  $v_i \in V$ . Thus,  $|l_i| = \sum_{j=1}^n (i, j)$  and  $|L| = \sum_{i=1}^n |l_i|$ . Mathematically noted,  $|l_i| = m_i$  is said to be the degree of  $v_i$ , while  $m = |L|/2$  is said to be the volume of  $G$ . Furthermore, the strength of each node  $i$  can be specified in more details by  $|l_i| = |l_{i,intra}|_{i \in C} + |l_{i,inter}|_{i \in C}$ , where  $|l_{i,intra}|_{i \in C}$  and  $|l_{i,inter}|_{i \in C}$  be the intra-strength and inter-strength of node  $i$ , respectively. Generalizing this to all nodes, implies  $|L| = |L_{intra}| + |L_{inter}| = \sum_{i=1}^n |l_{i,intra}|_{i \in C} + \sum_{i=1}^n |l_{i,inter}|_{i \in C}$ .

**3 Modularity based co-clustering model**

Co-clustering or simultaneous matrix partitioning (in contrast to clustering, as depicted in Figure-1) needs a quality function that can capture the embedded *distinct* sub-matrix structures. The *modularity* (normally noted as  $Q$ ) model defined after Newman and Girvan, lays the foundation of many existing successful graph clustering algorithms [5], [6]. The purpose of  $Q$  is to capture the hidden structure of sub-graphs or community sets in complex networks by maximizing intra-cluster links while minimizing inter-cluster ones.

Consider partitioning  $V$  of  $G$  into a co-clustering solution  $\mathcal{C} = \{C_1, \dots, C_K\}$  such that each vertex  $v_i, 1 \leq i \leq n$  is exactly assigned to one cluster  $C_j, 1 \leq j \leq K$ . The impact of  $E$  in  $\mathcal{C}$  is quantified in two distinct terms. The set of edges between vertices existing in two distinct clusters:  $E(C_i, C_j), 1 \leq i, j \leq K$  and  $i \neq j$  and the set of edges found inside one cluster:  $E(C_i, C_i), 1 \leq i \leq K$ . Then, modularity will award  $\mathcal{C}$  according to the fraction of connections inside its communities as formulated in Eq. 1. The left term in Eq. 1 biases towards a solution  $\mathcal{C}$  that is covered with a set of densely intra-connected modules, i.e. many edges fall within each sub-graph  $\{C_1, \dots, C_K\}$ . On the other hand, the right term in Eq. 1 expresses that the expected value of the same edge density in  $\mathcal{C}$  with the same community structure  $\{C_1, \dots, C_K\}$  but fall at random between the vertices should be small.  $Q$  will approach its minimum at 0 if the number of within-community edges is no better than random. On the other hand, values approaching  $Q = 1$ , which is the maximum, indicate strong community structure.

$$Q(\mathcal{C}) = \sum_{i=1}^K \left[ \frac{|E(C_i, C_i)|}{m(\mathcal{C})} - \left( \frac{\sum_{v \in C_i} m(v)}{2m(\mathcal{C})} \right)^2 \right] \tag{1}$$



**Figure 1**-Clustering against co-clustering. Left: clustering means partitioning all data vectors with all their features into  $k$  (sometimes unknown) disjoint groups. Right: co-clustering, or bi-clustering, means partitioning into a set of  $k$  (sometimes unknown) blocks each containing a consistent local feature pattern (Note that it is not generally possible to display several bi-clusters at the same time as contiguous blocks).

#### 4 The proposed heuristic based modularity

This section introduces a heuristic based approach for modularity with three different optimization models. A set of protein-neighborhood related functions is proposed to extend, accordingly, the unveil ability of a single objective EA. First, the main components that characterize the evolution process of EA (solution representation and perturbation operators) are formulated towards solving the problem. Then, the optimization models and the heuristic operator are introduced and formulated to improve the quality of generated complexes in the search space. Finally, the main steps of the proposed EA is outlined.

##### 4.1 The proposed EA

Any Evolutionary Algorithm (EA) searches for appropriate solutions from the set  $\Omega$  of all possible solutions of the problem at hand. Generally, the search for good solutions is performed through individual evaluations, selection, crossover, and mutation operators. The design of such operators would then determine the characteristic of the adopted EA. In this section, the definitions of all components are relaxed for the purpose of complex detection problem in PPI networks.

First, the construction of several, but unknown, number of complexes among the interacted proteins of a given PPI network, is an important issue that the individual representation (i.e. chromosome genotype encoding) should take quite seriously. In the proposed EA, the locus-based representation used in [7] is adopted. A chromosome  $P$  of the population  $\mathbb{P}$  is defined as a collection of node-node neighbor genes. A single gene in the chromosome  $P$  is defined by its locus and its allele. Consider a PPI network  $\mathcal{N}$  with  $n$  proteins. Then,  $P$  will consist of  $n$  genes, where locus  $i$  identifies protein  $i$  in the network, while its allele value  $j$  corresponds to the neighbor  $j$  that has an actual interaction with node  $i$  in the network, i.e.  $(v_i, v_j) \in E$ . This in turn implies that both proteins  $i$  and  $j$  will be in the same complex. The decoding function  $\delta$  of a chromosome  $P$  (chromosome phenotype) outlines one of the possible partitioning of the network  $\mathcal{N}$  into complexes, i.e.  $\delta(P): \mathcal{C} = \{\mathcal{C}\}_{i=1}^K$ . However,  $K$  could vary from one chromosome to another.

Once the population is created and their individuals are evaluated (according to the modularity in Eq. 1), a set of good population of parents is selected and processed by perturbation operators to create better child individuals. Two main perturbation operators are used. These are crossover  $\Psi_x$  and mutation  $\Psi_m$ .

Uniform crossover is used and achieved with a specified *chromosome-wise* crossover probability,  $p_x$ . Consider two chromosomes  $P_1: (I_{1,1}, I_{1,2}, \dots, I_{1,n})$  and  $P_2: (I_{2,1}, I_{2,2}, \dots, I_{2,n})$  to be the two participating parents in the crossover.

With probability  $p_x$ , a child  $P': (I'_1, I'_2, \dots, I'_n)$  can be generated from the two parents by uniformly mixing their alleles (i.e. performing *protein-wise fair* combination). This can be formally defined by:

$$\Psi_x: (P_1, P_2, p_x) \rightarrow P'$$

$$\forall i, 1 \leq i \leq n:$$

$$I'_i = \begin{cases} I_{1,i} & \text{if } r \leq 0.5 \\ I_{2,i} & \text{otherwise} \end{cases}$$

(2)

where  $r \sim [0,1]$  is a uniform random number.

The mutation operator  $\Psi_m$  imitates the traditional *allelic* mutation operator which works on allele values and alters, with a specified mutation probability  $p_m$ , the allele (i.e. neighbor) of a selected locus (i.e. node). This can formally be specified by:  $\Psi_m: (P, p_m) \rightarrow P'$

$$\forall i, 1 \leq i \leq n \wedge r \leq p_m:$$

$$I'_i = j | (i, j) \in E$$

(3)

where  $r \sim [0,1]$  is a uniform random number.

##### 4.2 Formulation of the heuristic optimization models and operator

Generally, heuristic operator or search heuristic is defined to be a rule that decides which solution, given the current solution, to generate or to visit next based on some heuristic criterion. In evolutionary computation community, the search for designing appropriate heuristic for a given problem is essential and can harness the performance of the algorithm. In the following discussion, we introduce a heuristic operator with three optimization models tailored, here, specifically for complex detection problem.

The general characteristic of complexes in PPI networks expresses dense interactions within complexes while more sparse interactions among different complexes. The main purpose of the proposed heuristic operator  $\Psi_h: P \rightarrow P'$  is to move proteins between the complexes of an individual solution  $P$ . The movement of the selected proteins should reduce the problems of both sparse intra-connections and dense inter-ties. Thus, the proposed  $\Psi_h$  operator works, with a specified probability  $p_h$ , on those nodes maintaining un-reliable template in their chosen complexes and move them to other complexes that can participate within their proteins more reliably. Here, we propose three different optimization models on how to define reliability assignment of a node to a given complex.

Consider an individual chromosome  $P: (I_1, I_2, \dots, I_n)$  corresponding to a candidate partitioning solution  $\mathcal{C} = \{C_1, \dots, C_K\}$  with  $K$  complexes. Let node  $i$  that corresponds to gene  $I_i$  being located in complex  $k$ , where  $1 \leq k \leq K$ . Then, node  $i$  inside complex  $k$  has a possible reliable interaction assignment that can be expressed as the difference between the impact of the intra-connections and inter-connections:

$$\text{ReliableAssign}(i, C_k) = \text{IntraImpact}_{i \in C_k} - \text{InterImpact}_{i \in C_k} \quad (4)$$

For *IntraImpact* and *InterImpact* in Eq. 4, three models are proposed to define them. Let us first consider an extended version of the adjacency matrix  $A$  (discussed in Section 2). A weighted adjacency matrix  $wA$  is constructed using Eq. 5.

$$wA(i, j) = \begin{cases} \sum_{k=1}^n \left( \frac{A(i, k)}{|i|-1} \right) \times \left( \frac{A(k, j)}{|k|} \right) & \text{if } A(i, j) = 1 \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

The proposed heuristic operator,  $\Psi_h$ , (see Algorithm 1), then, moves node  $i$  to another complex  $k'$ ,  $1 \leq k' \leq K$  and  $k' \neq k$  where node  $i$  could maintain, there, the highest reliability assignment, i.e. with the highest difference between intra-connections and inter-connections impact. The proposed heuristic operator ( $\Psi_h$ ), then, can be stated formally as in Eq. 5. Note that when more than one complex can receive node  $i$  with equal *ReliableAssign* value, then  $\Psi_h$  randomly selects any one of these complexes.

$$\Psi_h(i \in C_k, p_h) = \max_{C_{k'} \in \mathcal{C}} \text{ReliableAssign}(i, C_{k'}) \quad (6)$$

Now, three different heuristic models are proposed to reflect *IntraImpact* <sub>$i \in C_k$</sub>  and *InterImpact* <sub>$i \in C_k$</sub>  terms of *ReliableAssign*. The first model (*heuristic#1*) considers the difference between the accumulated impact of the weighted intra-connections of all proteins  $i \in C_k$  and their inter-connections, i.e.:

$$\text{IntraImpact}_{i \in C_k} = \sum_{i \in C_k, j \in C_k} wA(i, j) \quad (7)$$

$$\text{InterImpact}_{i \in C_k} = \sum_{i \in C_k, j \notin C_k} wA(i, j) \quad (8)$$

Equation 7 and 8 are also maintained in the next two heuristic models, however, to be combined with additional terms in an alternating fashion with equal probability. For the second heuristic model (*heuristic#2*), *IntraImpact* <sub>$i \in C_k$</sub>  and *InterImpact* <sub>$i \in C_k$</sub> , are computed with equal probability, to be either: 1) Eq.7 and Eq.8, respectively, or 2) to return accumulated impact of the intra-connections of neighbors of all proteins  $i \in C_k$  and their inter-connections, i.e.:

$$\text{IntraImpact}_{i \in C_k} = \sum_{i \in C_k, j \in C_k \wedge (i, j) \in E} A(i, j) \quad (9)$$

$$\text{InterImpact}_{i \in C_k} = \sum_{i \in C_k, j \notin C_k \wedge (i, j) \in E} A(i, j) \quad (10)$$

Finally, the third model (*heuristic#3*) computes *IntraImpact* <sub>$i \in C_k$</sub>  and *InterImpact* <sub>$i \in C_k$</sub> , with equal probability, to be either: 1) equal to Eq.7 and Eq.8, respectively, or 2) to reflect the accumulated impact of the intra-connections and inter-connections of all proteins  $i \in C_k$ , i.e.:

$$\text{IntraImpact}_{i \in C_k} = \sum_{i \in C_k, j \in C_k} A(i, j) \quad (11)$$

$$\text{InterImpact}_{i \in C_k} = \sum_{i \in C_k, j \notin C_k} A(i, j) \quad (12)$$

**Algorithm 1:  $\Psi_h$** 

**Input:** 1) chromosome phenotype  $\mathcal{C} = \{C_1, \dots, C_K\}$ ,  
 2) number of proteins  $n$ ,  
 3) probability of heuristic operator  $p_h$

**Output:** modified chromosome phenotype  $\mathcal{C}' = \{C_1, \dots, C_K\}$

```

1  for  $i = 1$  to  $n$  do
2      if ( $rand \leq p_h$ ) // apply heuristic movement to protein  $i$ 
3          set  $C_i \leftarrow Complex(i)$ ; // return the current complex of protein  $i$ 
4          compute  $cur\_intra\_impact\_i$ ; // according to the heuristic model
5          compute  $cur\_inter\_impact\_i$ ; // according to the heuristic model
6          set  $[C_j, ReliableAssign] \leftarrow argmax_{C_j \in \mathcal{C}} (ReliableAssign(i, C_j))$ ; //return complex
            $C_j$  with maximum reliability assignment for protein  $i$ 
7          if ( $ReliableAssign > (cur\_intra\_impact\_i - cur\_inter\_impact\_i)$ )
8              set  $Complex(i) \leftarrow C_j$ ;
9          end if
10     end if
11 end for

```

**4.3 General EA layout**

The overall component of the proposed EA with the proposed heuristic model is then presented in Algorithm 2.

**Algorithm 2: heuristic EA for complex detection problem in PPI networks**

**Input:** 1) PPI network:  $\mathcal{N}(n, E)$ ,  
 2) population size:  $\mu$ , and maximum number of generations  $max_t$   
 3) EA operators and their probabilities:  $s, \Psi_x, \Psi_m, \Psi_h, p_x, p_m, p_h$ ,

**Output:** Best individual solution  $P$

```

1  initialize  $\mathbb{P} \leftarrow \{P_1, P_2, \dots, P_\mu\}$ ;
2   $t \leftarrow 0$ ;
3  evaluate modularity for each individual in the population  $\mathbb{P}(t) = \{Q_1, Q_2, \dots, Q_\mu\}$ ;
4  while ( $t \leq max_t$ ) do
5      for  $i \leftarrow 1$  to  $\mu$  do
6           $P_{i,1}(t) \leftarrow select(\mathbb{P}_i(t))$ ; // select parent 1 for  $P_i$ 
7           $P_{i,2}(t) \leftarrow select(\mathbb{P}_i(t))$ ; // select parent 2 for  $P_i$ 
8           $P_i(t) \leftarrow \Psi_x(P_{i,1}(t), P_{i,2}(t), p_x)$ ;
9           $P_i(t) \leftarrow \Psi_m(P_i(t), p_m)$ ;
10          $P_i(t) \leftarrow \Psi_h(P_i(t), p_h)$ ;
11         evaluate  $Q(P_i(t))$ ;
12     end for
13      $t \leftarrow t + 1$ ;
14 end while
15 return  $P \in \mathbb{P}(t)$  with maximum  $Q$ ;

```

### 5 Results

The experiments include two commonly PPI networks, denoted by  $\mathcal{N}1$  and  $\mathcal{N}2$ .  $\mathcal{N}1$  has  $n = 990$  proteins with  $E = 4687$  interactions, while  $\mathcal{N}2$  has  $n = 1443$  proteins and  $E = 6993$  interactions. To validate the quality of the predicted complexes generated by the tested EA without heuristic against with heuristic, two sets of golden standard complexes ( $Cmplx\_D1$  and  $Cmplx\_D2$ ) drawn from the Munich Information Center for Protein Sequence (MIPS) catalog [8] are used in the experiments.  $Cmplx\_D1$  contains 81 complexes, while  $Cmplx\_D2$  is made of 162 hand-curated complexes To evaluate the quality of the detected complexes obtained by the EA, several metrics are used. The predicted set of complexes  $\mathcal{C} = \{C_1, C_2, \dots, C_K\}$  obtained by EA is compared with the golden standard complexes  $\mathcal{C}^* = \{C_1^*, C_2^*, \dots, C_{K^*}^*\}$  of  $K^*$  complexes. A predicted complex  $C_i$  in the solution  $\mathcal{C}$  overlaps a golden standard complex  $C_j^*$  by an overlapping score ( $OS$ ). Then, the predicted complex  $C_i$  matches the golden standard complex  $C_j^*$  if  $OS$  is equal or larger than a specified threshold,  $\sigma_{OS}$ , [9].

$$OS(C_i, C_j^*) = \frac{|C_i \cap C_j^*|^2}{|C_i| |C_j^*|} \tag{13}$$

where  $|C_i \cap C_j^*|$  is the number of proteins common to both a predicted complex and a golden standard complex.

$$match(C_i, C_j^*) = \begin{cases} 1 & \text{if } OS(C_i, C_j^*) \geq \sigma_{OS} \\ 0 & \text{otherwise} \end{cases} \tag{14}$$

At both complex and protein levels, the three standard metrics of recall, precision, and F measure are evaluated. The complex/protein levels,  $recall/recall_N$ ,  $precision/precision_N$ , and cumulative  $F\text{-measure}/F_N$  are defined. In  $recall$ , the fraction of golden complexes/proteins that are matched to any predicted complex is determined. On the other hand,  $precision$  refers to the fraction of predicted complexes/proteins that are matched to any golden standard complex. A harmonic mean of both  $recall$ , and  $precision$  is reflected by  $F\text{-measure}$ .

$$recall = \frac{|C_i^* : C_i^* \in \mathcal{C}^* \wedge \exists C_j \in \mathcal{C} \rightarrow match(C_i^*, C_j)|}{K^*} \tag{15}$$

$$precision = \frac{|C_i : C_i \in \mathcal{C} \wedge \exists C_j^* \in \mathcal{C}^* \rightarrow match(C_i, C_j^*)|}{K_C} \tag{16}$$

$$F\text{-measure} = \frac{2 * recall * precision}{recall + precision} \tag{17}$$

$$recall_N = \frac{\sum_{i=1}^{K^*} |m_i|}{\sum_{i=1}^{K^*} |C_i^*|} \tag{18}$$

where  $|m_i| = \max_{C_j \in \mathcal{C}} |match(C_i^*, C_j)|$ .

$$precision_N = \frac{\sum_{i=1}^{K_C} |m_i|}{\sum_{i=1}^{K_C} |C_i|} \tag{19}$$

where  $|m_i| = \max_{C_j^* \in \mathcal{C}^*} |match(C_i, C_j^*)|$ .

$$F_N = \frac{2 * recall_N * precision_N}{recall_N + precision_N} \tag{20}$$

Other measures can be computed with no dependency to the overlapping score ( $\sigma_{OS}$ ). These are general sensitivity ( $sensitivity$ ), general positive predictive value ( $PPV$ ), and  $accuracy$  [9]. General sensitivity ( $sensitivity$ ) between the set of complexes  $\mathcal{C}^* = \{C_1^*, C_2^*, \dots, C_{K^*}^*\}$  and the set of detected partitioning solution  $\mathcal{C} = \{C_1, C_2, \dots, C_K\}$  is the weighted average of complex-wise sensitivity of all reference complexes (Eq. 21). Similarly, general  $PPV$ , with respect to the detected complexes (Eq. 22).

$$sensitivity = \frac{\sum_{i=1}^{K^*} \max_{j=1}^K \frac{T(i,j)}{n_i^*}}{\sum_{i=1}^{K^*} n_i^*} \tag{21}$$

$$PPV = \frac{\sum_{j=1}^K \max_{i=1}^{K^*} \frac{T(i,j)}{T_j}}{\sum_{j=1}^K \sum_{i=1}^{K^*} T(i,j)} \tag{22}$$

where  $T_j$  represents the marginal sum of column  $j$ . The tradeoff between *sensitivity* and *PPV* can be represented by the geometric *accuracy*. High accuracy (Eq. 23) value requires a high performance for both *sensitivity* and *PPV*.

$$accuracy = \sqrt{sensitivity * PPV} \tag{23}$$

Results for all mentioned metrics are reported in Tables-(1 – 7) and Figures-(2, 3). The reported results (given in bold) clearly reveal the positive impact of the proposed heuristic operator with the three different versions of models. The proposed heuristic operator extends the applicability of the well-known modularity function ( $Q$ ) to partition a given PPI network.

**Table 1-**Performance in terms of *recall* for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

$\sigma_{OS}$	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #1	heuristic #2	heuristic #3
0.1	0.8500	<b>0.8949</b>	<b>0.8538</b>	<b>0.8744</b>	0.9287	0.9267	0.9167	0.9193
0.15	0.7910	<b>0.8359</b>	0.7795	<b>0.8141</b>	0.8493	<b>0.8493</b>	0.8307	0.8420
0.2	0.7244	<b>0.7872</b>	<b>0.7256</b>	<b>0.7538</b>	0.7860	<b>0.7860</b>	0.7613	0.7653
0.25	0.6872	<b>0.7449</b>	<b>0.6782</b>	<b>0.7077</b>	0.7260	<b>0.7307</b>	0.7013	0.7147
0.3	0.6462	<b>0.7090</b>	<b>0.6551</b>	<b>0.6731</b>	0.6580	0.6573	0.6307	0.6467
0.35	0.6269	<b>0.7038</b>	<b>0.6359</b>	<b>0.6603</b>	0.5767	<b>0.5813</b>	0.5660	<b>0.5787</b>
0.4	0.5936	<b>0.6872</b>	<b>0.6064</b>	<b>0.6244</b>	0.5333	<b>0.5367</b>	<b>0.5340</b>	<b>0.5413</b>
0.45	0.5577	<b>0.6538</b>	<b>0.5808</b>	<b>0.6000</b>	0.4633	<b>0.4673</b>	<b>0.4660</b>	<b>0.4773</b>
0.5	0.5423	<b>0.6346</b>	<b>0.5603</b>	<b>0.5872</b>	0.4440	<b>0.4467</b>	<b>0.4467</b>	<b>0.4520</b>

**Table 2-**Performance in terms of *recall<sub>N</sub>* for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

$\sigma_{OS}$	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #1	heuristic #2	heuristic #3
0.1	0.8582	<b>0.8936</b>	<b>0.8835</b>	<b>0.8970</b>	0.6062	<b>0.6070</b>	<b>0.6121</b>	<b>0.6116</b>
0.15	0.8180	<b>0.8498</b>	<b>0.8300</b>	<b>0.8505</b>	0.5833	0.5790	0.5800	<b>0.5856</b>
0.2	0.7575	<b>0.8137</b>	<b>0.7866</b>	<b>0.8016</b>	0.5571	0.5507	0.5459	0.5518
0.25	0.6983	<b>0.7635</b>	<b>0.7100</b>	<b>0.7466</b>	0.5220	0.5155	0.5101	0.5126
0.3	0.6163	<b>0.7091</b>	<b>0.6779</b>	<b>0.6765</b>	0.4717	0.4707	0.4542	0.4650
0.35	0.6011	<b>0.7049</b>	<b>0.6535</b>	<b>0.6668</b>	0.4026	<b>0.4113</b>	<b>0.4138</b>	<b>0.4301</b>
0.4	0.5478	<b>0.6860</b>	<b>0.6029</b>	<b>0.6157</b>	0.3576	0.3740	<b>0.3783</b>	<b>0.3916</b>
0.45	0.4964	<b>0.6288</b>	<b>0.5524</b>	<b>0.5601</b>	0.3040	<b>0.3211</b>	<b>0.3215</b>	<b>0.3403</b>
0.5	0.4718	<b>0.5845</b>	<b>0.5180</b>	<b>0.5382</b>	0.2783	<b>0.2881</b>	<b>0.2904</b>	<b>0.3042</b>



**Table 3** Performance in terms of *precision* for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

$\sigma_{OS}$	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #	heuristic #	heuristic #
0.1	0.7735	0.7588	<b>0.7859</b>	<b>0.7799</b>	0.5932	<b>0.5964</b>	<b>0.6125</b>	<b>0.6049</b>
0.15	0.7361	0.7212	<b>0.7533</b>	<b>0.7431</b>	0.5656	<b>0.5698</b>	<b>0.5897</b>	<b>0.5872</b>
0.2	0.7332	0.7175	<b>0.7533</b>	<b>0.7402</b>	0.5403	<b>0.5458</b>	<b>0.5643</b>	<b>0.5634</b>
0.25	0.7226	0.7097	<b>0.7485</b>	<b>0.7343</b>	0.4964	0.4945	<b>0.5169</b>	<b>0.5235</b>
0.3	0.7002	0.6931	<b>0.7421</b>	<b>0.7224</b>	0.4913	0.4814	<b>0.5068</b>	<b>0.5199</b>
0.35	0.6866	<b>0.6880</b>	<b>0.7343</b>	<b>0.7163</b>	0.4694	0.4548	<b>0.4878</b>	<b>0.4924</b>
0.4	0.6718	<b>0.6844</b>	<b>0.7172</b>	<b>0.7059</b>	0.4559	0.4331	<b>0.4751</b>	<b>0.4764</b>
0.45	0.6452	<b>0.6653</b>	<b>0.7016</b>	<b>0.6869</b>	0.4178	0.4039	<b>0.4441</b>	<b>0.4428</b>
0.5	0.6288	<b>0.6477</b>	<b>0.6802</b>	<b>0.6735</b>	0.4103	0.3956	<b>0.4378</b>	<b>0.4339</b>

**Table 4**-Performance in terms of *precision<sub>N</sub>* for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

$\sigma_{OS}$	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #	heuristic #	heuristic #
0.1	0.6338	<b>0.7185</b>	<b>0.6683</b>	<b>0.6779</b>	0.6829	<b>0.6976</b>	<b>0.6890</b>	<b>0.6974</b>
0.15	0.6302	<b>0.7138</b>	<b>0.6653</b>	<b>0.6743</b>	0.6768	<b>0.6929</b>	<b>0.6853</b>	<b>0.6944</b>
0.2	0.6300	<b>0.7107</b>	<b>0.6653</b>	<b>0.6715</b>	0.6672	<b>0.6841</b>	<b>0.6735</b>	<b>0.6829</b>
0.25	0.6097	<b>0.7037</b>	<b>0.6511</b>	<b>0.6653</b>	0.6510	<b>0.6607</b>	0.6509	<b>0.6687</b>
0.3	0.5692	<b>0.6779</b>	<b>0.6359</b>	<b>0.6325</b>	0.6429	0.6328	0.6241	<b>0.6465</b>
0.35	0.5596	<b>0.6743</b>	<b>0.6194</b>	<b>0.6260</b>	0.6106	0.6082	<b>0.6146</b>	<b>0.6246</b>
0.4	0.5297	<b>0.6719</b>	<b>0.5867</b>	<b>0.6067</b>	0.5826	0.5710	<b>0.5910</b>	<b>0.5894</b>
0.45	0.4906	<b>0.6282</b>	<b>0.5501</b>	<b>0.5588</b>	0.5259	<b>0.5312</b>	<b>0.5331</b>	<b>0.5510</b>
0.5	0.4688	<b>0.5845</b>	<b>0.5180</b>	<b>0.5382</b>	0.5078	<b>0.5100</b>	<b>0.5138</b>	<b>0.5294</b>

**Table 5-**Performance in terms of *Fmeasure* for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

$\sigma_{OS}$	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #	heuristic #	heuristic #
0.1	0.8095	<b>0.8209</b>	<b>0.8178</b>	<b>0.8241</b>	0.7238	<b>0.7255</b>	<b>0.7342</b>	<b>0.7296</b>
0.15	0.7621	<b>0.7737</b>	<b>0.7655</b>	<b>0.7765</b>	0.6787	<b>0.6818</b>	<b>0.6895</b>	<b>0.6918</b>
0.2	0.7281	<b>0.7498</b>	<b>0.7385</b>	<b>0.7467</b>	0.6401	<b>0.6440</b>	<b>0.6478</b>	<b>0.6489</b>
0.25	0.7041	<b>0.7260</b>	<b>0.7110</b>	<b>0.7204</b>	0.5894	<b>0.5895</b>	<b>0.5948</b>	<b>0.6042</b>
0.3	0.6716	<b>0.7002</b>	<b>0.6954</b>	<b>0.6965</b>	0.5624	0.5554	0.5615	<b>0.5762</b>
0.35	0.6550	<b>0.6950</b>	<b>0.6811</b>	<b>0.6866</b>	0.5172	0.5098	<b>0.5234</b>	<b>0.5319</b>
0.4	0.6298	<b>0.6849</b>	<b>0.6568</b>	<b>0.6622</b>	0.4912	0.4789	<b>0.5022</b>	<b>0.5066</b>
0.45	0.5980	<b>0.6587</b>	<b>0.6351</b>	<b>0.6401</b>	0.4391	0.4329	<b>0.4543</b>	<b>0.4592</b>
0.5	0.5821	<b>0.6400</b>	<b>0.6140</b>	<b>0.6269</b>	0.4262	0.4191	<b>0.4418</b>	<b>0.4426</b>

**Table 6-**Performance in terms of  $F_N$  for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

$\sigma_{OS}$	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #	heuristic #	heuristic #
0.1	0.7286	<b>0.7962</b>	<b>0.7609</b>	<b>0.7721</b>	0.6422	<b>0.6489</b>	<b>0.6480</b>	<b>0.6516</b>
0.15	0.7115	<b>0.7757</b>	<b>0.7383</b>	<b>0.7517</b>	0.6264	<b>0.6305</b>	<b>0.6281</b>	<b>0.6353</b>
0.2	0.6873	<b>0.7584</b>	<b>0.7205</b>	<b>0.7302</b>	0.6071	<b>0.6099</b>	0.6027	<b>0.6104</b>
0.25	0.6506	<b>0.7322</b>	<b>0.6792</b>	<b>0.7031</b>	0.5791	0.5789	0.5716	<b>0.5802</b>
0.3	0.5917	<b>0.6930</b>	<b>0.6562</b>	<b>0.6536</b>	0.5438	0.5396	0.5255	0.5408
0.35	0.5796	<b>0.6891</b>	<b>0.6359</b>	<b>0.6456</b>	0.4847	<b>0.4904</b>	<b>0.4940</b>	<b>0.5093</b>
0.4	0.5385	<b>0.6788</b>	<b>0.5946</b>	<b>0.6111</b>	0.4426	<b>0.4517</b>	<b>0.4607</b>	<b>0.4704</b>
0.45	0.4934	<b>0.6285</b>	<b>0.5512</b>	<b>0.5594</b>	0.3848	<b>0.3999</b>	<b>0.4005</b>	<b>0.4206</b>
0.5	0.4702	<b>0.5845</b>	<b>0.5180</b>	<b>0.5382</b>	0.3591	<b>0.3677</b>	<b>0.3705</b>	<b>0.3860</b>

**Table 7** Performance in terms of *Sensitivity*, *PPV*, *Accuracy*, *CCF* and *Strength* for  $\mathcal{N}1$  and  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics.

	$\mathcal{N}1$				$\mathcal{N}2$			
	Q without heuristic	Q with			Q without heuristic	Q with		
		heuristic #1	heuristic #2	heuristic #3		heuristic #1	heuristic #2	heuristic #3
<i>Sensitivity</i>	0.9621	0.9613	<b>0.9772</b>	<b>0.9752</b>	0.6420	0.6411	<b>0.6531</b>	<b>0.6527</b>
<i>PPV</i>	0.6345	<b>0.7201</b>	<b>0.6688</b>	<b>0.6790</b>	0.2785	<b>0.2810</b>	0.2778	<b>0.2832</b>
<i>Accuracy</i>	0.7810	<b>0.8317</b>	<b>0.8081</b>	<b>0.8134</b>	0.4229	<b>0.4244</b>	<b>0.4259</b>	<b>0.4299</b>
<i>CCF</i>	559.6000	<b>589.3500</b>	<b>576.9000</b>	<b>579.8000</b>	875.4500	<b>876.9000</b>	<b>885.2500</b>	<b>890.1000</b>
<i>Strength</i>	0.6647	0.6635	<b>0.7435</b>	<b>0.7312</b>	0.5647	0.5636	<b>0.6098</b>	<b>0.6311</b>

Two additional metrics are also included in Table 7. These are Cross Common Fraction (*CCF*) and Strength of complex structure. *CCF* compares each pair of complexes, in which one comes from the golden data ( $C_i$ ) and the second comes from the detected result ( $C_j$ ), to find the maximal shared parts. *Strength* measures the intensity of the detected complexes ( $C_j$ ). This measure comes after Radicchi et al. [10] well-known community definition. They showed that a community structure usually seems correct, strong, and valid if most members with their neighbors are inside one community.

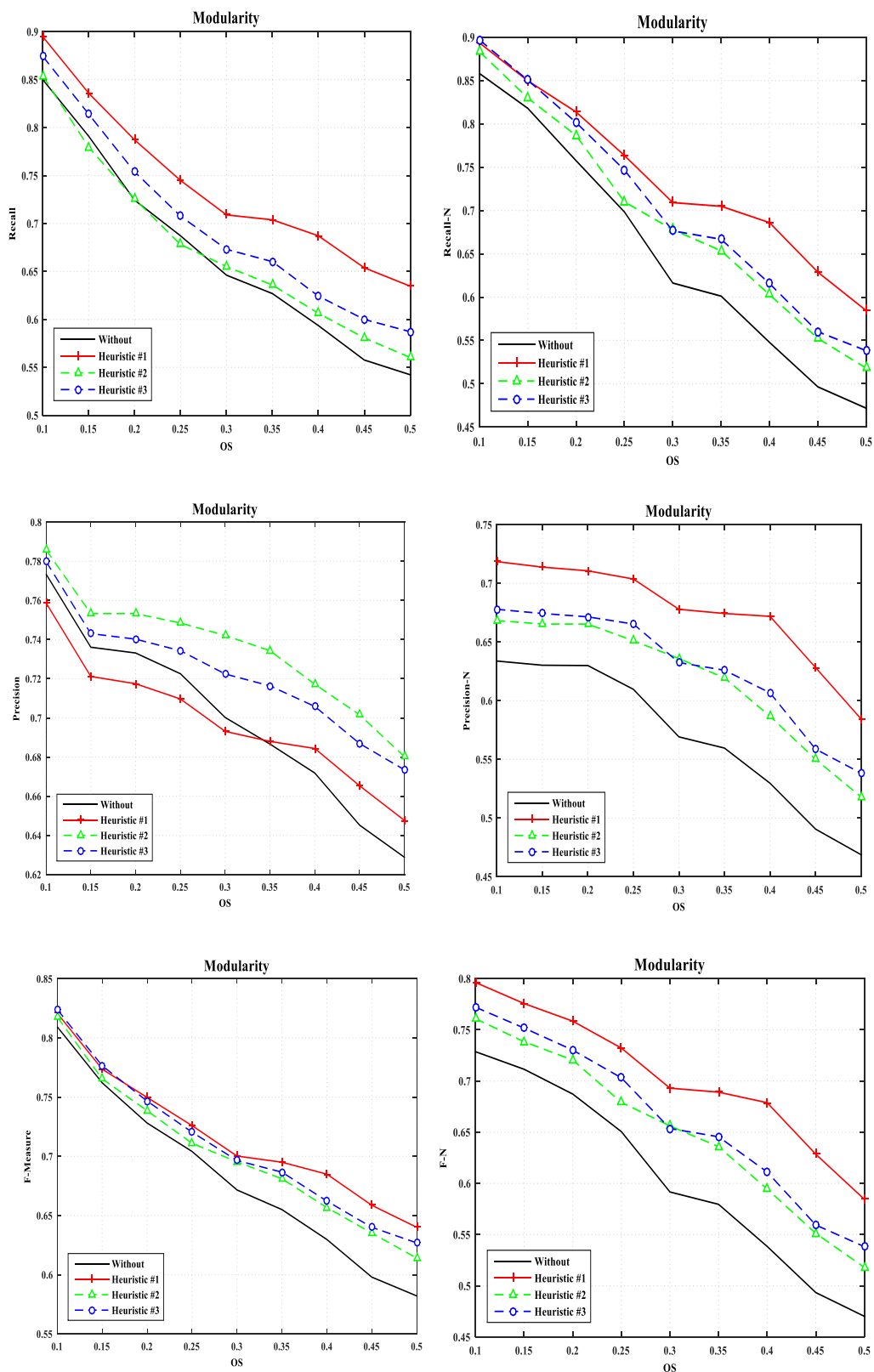
$$CCF = \frac{1}{2} \sum_{i=1}^{K^*} \max_j |C_i \cap C_j| + \frac{1}{2} \sum_{j=1}^K \max_i |C_i \cap C_j| \tag{24}$$

$$Strength = \frac{1}{K} \sum_{i=1}^K Score(C_i) \tag{25}$$

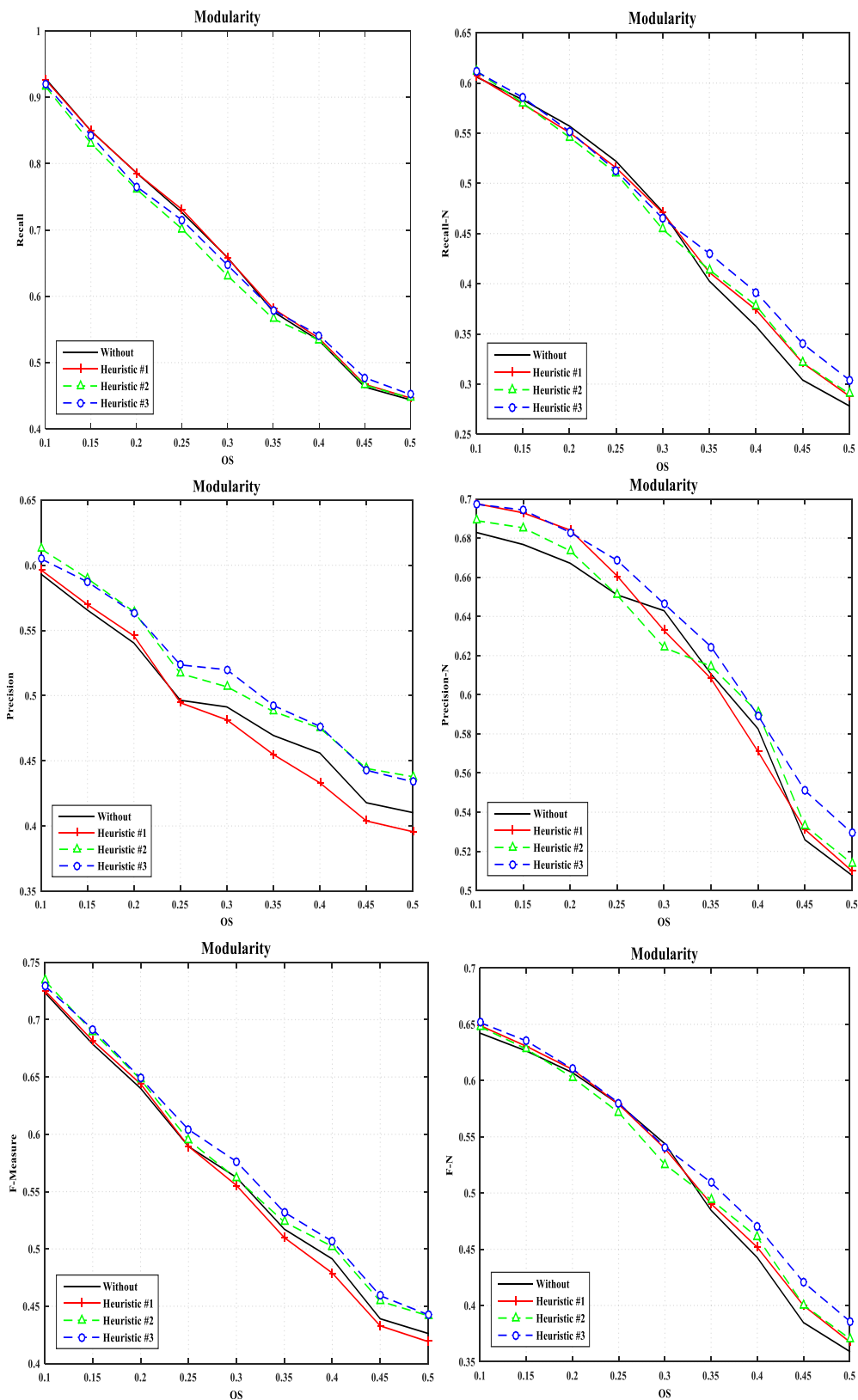
where:

$$Score(C_i) = \begin{cases} 1 & \text{if } \forall v \in C_i \rightarrow |l_{v,intra}|_{v \in C_i} + |l_{v,inter}|_{v \in C_i} \\ 0.5 & \text{if } \sum_{\forall v \in C_i} |l_{v,intra}|_{v \in C_i} > \sum_{\forall v \in C_i} |l_{v,inter}|_{v \in C_i} \\ 0 & \text{otherwise} \end{cases} \tag{26}$$

Generally, all reported results reveal consistent effectiveness of the proposed heuristic models and the proposed operator when coupled with the modularity function to partition a PPI network into different complexes. The detected complexes have satisfactory aggregation of intra connections where more intra-group connections appear than inter-group connections. This occurs much clearly in the first PPI network than in the second network. This is mainly due to the fact that the second network contains several overlapping proteins.



**Figure 2-**Performance at the complex and protein levels for  $\mathcal{N}1$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.



**Figure 3-**Performance at the complex and protein levels for  $\mathcal{N}2$  when  $Q$  is adopted without heuristic, in one hand, and with the proposed heuristics, on the other hand. Threshold of overlapping score  $\sigma_{OS}$  is varied from 0.1 to 0.5 in step of 0.05.

## 6 Conclusions

The results reported in this paper show the importance on explicitly considering neighboring relations and heuristic movement models to improve the detection reliability of the modularity based EA. The proposed heuristic models emphasize the complex oriented structures where dense intra connections and sparse inter connections are declared. Generally, the results provided in this paper are encouraging, especially for the future of the proposed heuristic operators and models, and their extensions in the application to other complex detection models (e.g. community score, normalized cut, and ratio cut) and difficult mining problems in complex networks, for example for community detection problem in weighted networks and overlapping complex detection.

## References

1. Lancichinetti, A. and Fortunato, S. **2009**. Community detection algorithms: a comparative analysis. *Physical review E*, **80**(5): 056117.
2. Brandes, U., Delling, D., Gaertler, M., Gorke, R., Hoefer, M., Nikoloski, Z. and Wagner, D. **2008**. On modularity clustering. *IEEE transactions on knowledge and data engineering*, **20**(2): 172-188.
3. Clauset, A., Newman, M.E. and Moore, C. **2004**. Finding community structure in very large networks. *Physical review E*, **70**(6): 066111.
4. Coello, C.A.C., Van Veldhuizen, D.A. and Lamont, G.B. **2002**. *Evolutionary algorithms for solving multi – objective problems* (Vol. 242). New York: Kluwer Academic.
5. Newman, M.E.J and Girvan, M. **2004**. “Finding and evaluating community structure in networks,” *Physical Review E*, **69**(026113).
6. Girvan, M. and Newman, M.E. **2002**. Community structure in social and biological networks. *Proceedings of the national academy of sciences*, **99**(12): 7821-7826.
7. Attea, B.A., Hariz, W.A. and Abdulhalim, M.F. **2016**. Improving the performance of evolutionary multi-objective co-clustering models for community detection in complex social networks. *Swarm and Evolutionary Computation*, **26**: 137-156.
8. Gavin, A. C., Aloy, P., Grandi, P., Krause, R., Boesche, M., Marzioch, M. and Edlmann, A. 2006. Proteome survey reveals modularity of the yeast cell machinery. *Nature*, **440**(7084): 631-636.
9. Mewes, H. W., Frishman, D., Güldener, U., Mannhaupt, G., Mayer, K., Mokrejs, M. and Weil, B. **2002**. MIPS: a database for genomes and protein sequences. *Nucleic acids research*, **30**(1): 31-34.
10. Radicchi, F., Castellano, C., Cecconi, F., Loreto, V. and Parisi, D. **2004**. Defining and identifying communities in networks. *Proceedings of the National Academy of Sciences of the United States of America*, **101**(9): 2658-2663.