

## RESEARCH ARTICLE

# Robust Clustering for Ad Hoc Cognitive Radio Network

Di Li<sup>1\*</sup>, Erwin Fang<sup>2</sup>, James Gross<sup>3</sup>RWTH Aachen University<sup>1</sup>, Swisscom (Schweiz) AG<sup>2</sup>, KTH Royal Institute of Technology<sup>3</sup>

## ABSTRACT

Cluster structure in cognitive radio networks facilitates cooperative spectrum sensing, routing and other functionalities. Unlicensed channels, which are available for a group of cognitive radio users, consolidate the group into a cluster and the number of the available unlicensed channels decides that cluster's robustness against the licensed users' influence. This paper analyses the problem of how to form robust clusters in a cognitive radio network so that more cognitive radio users can get the benefits from cluster structure when the primary users' operations becomes more intense. We give a formal description of the robust clustering problem, prove it to be NP-hard and propose both centralized and distributed solutions. The congestion game model is adopted to analyze the process of cluster formation, which not only contributes to the design of the distributed clustering scheme, but also provides the guarantee on the convergence into Nash Equilibrium and the convergence speed. The proposed distributed clustering scheme outperforms state-of-the-art related works in terms of cluster robustness, convergence speed and overhead. Extensive simulations have been conducted which clearly support our claims. Copyright © 2017 John Wiley & Sons, Ltd.

### \* Correspondence

Chair of Communication and Distributed Systems Ahornstrae 55 - building E3 52074 Aachen Germany

Email: li@umic.rwth-aachen.de

## 1. INTRODUCTION

Cognitive radio (CR) is a promising technology to solve the spectrum scarcity problem for the upcoming era of internet of things [?, ?]. Licensed users access the spectrum allocated to them whenever there is information to be transmitted. In contrast, as one way, unlicensed users can access the spectrum via opportunistic spectrum access, i.e., they access the licensed spectrum only after validating the channel is unoccupied by licensed users, where spectrum sensing [?] plays an important role in this process. In this hierarchical spectrum access model [?], the licensed users are also called primary users (PU), while the unlicensed users are referred to as secondary users and constitute a so called cognitive radio network (CRN). Regarding the operation of CRN, efficient spectrum sensing is identified to be critical for a smooth operation of a cognitive radio network [?]. Efficient spectrum sensing can be achieved. The rate of false negative, i.e., misdetecting the active primary users, can be decreased by cooperative spectrum sensing of multiple secondary users, which has been shown to cope effectively with noise uncertainty and channel fading, thus remarkably improving the sensing accuracy [?]. cooperative Sensing <sup>Akyildiz11 \* 0. Collaborative sensing relies on the consensus of CR users\* within a certain area, in</sup>

the these selected cluster head. The size of  $C$  is denoted by  $|C|$ . When the cluster head of a cluster is  $i$ , we denote that cluster by  $C(i)$ .  $K(C)$  denotes the set of CCs in cluster  $C$ ,  $K(C) = \bigcap_{i \in C} K_i$ . The notations used in the system model are listed in Table ??.

### 2.1. Robust Clustering Problem in CRN

As introduced in Section ??, in order to be robust against primary users' activity, the formed clusters should have more CCs. On the other hand, the sizes of the formed clusters should be regulated, i.e., they don't diverge from a given value greatly.

**Definition 1.** Robust clustering problem in CRN.

As to a cognitive radio network where the set of CR nodes is  $N$ , the robust clustering problem is to decide the set of clusters  $\mathcal{T}$ , where

1. the intersection of any two clusters in  $\mathcal{T}$  is an empty set

**Table I.** Notations

Symbol	Description
$\mathcal{N}$	set of CR users in a CRN
$N$	number of CR users in a CRN, $N =  \mathcal{N} $
$\mathcal{K}$	set of licensed channels
$k(i)$	the working channel of user $i$
$\text{Nb}(i)$	the neighborhood of CR node $i$
$C(i)$	a cluster whose cluster head is $i$
$K_i$	the set of available channels at CR node $i$
$K(C(i))$	the set of available CCs of cluster $C(i)$
$h(C)$	the cluster head of a cluster $C$
$\delta$	the cluster size which is preferred
$S_i$	a set of claiming clusters, each of which includes debatable node $i$ after phase I
$d_i$	individual connectivity degree of CR node $i$
$g_i$	neighborhood connectivity degree of CR node $i$
$f(C)$	the number of CCs of a cluster $C$ , which is used in the problem description
$\mathcal{S}$	the collection of all the possible clusters in $\mathcal{N}$
$C_i$	the $i$ -th cluster in $\mathcal{S}$
$ C_i $	size of the cluster $C_i$
$ K(C_i) $	the number of CCs of cluster $C_i$
$n$	the number of debatable nodes
$m$	the number of claiming cluster heads

2. the union of clusters in  $\mathcal{T}$  is  $\mathcal{N}$

3. when the number of common channels for cluster  $C$  is denoted as  $f(C)$ , the sum of the  $f(C)$  is the maximal, where  $C \in \mathcal{T}$  and meanwhile the cluster sizes fall in the scope  $[\delta_1, \delta_2]$ .  $\delta_1, \delta_2 \in \mathbb{Z}^+$  and  $\delta_1 \leq \delta_2$ . When the cluster size is out of  $[\delta_1, \delta_2]$ ,  $f(C)$  is defined as 0.

4. the size of  $C$  in  $\mathcal{T}$  is allowed to be 1.

The decision version of this problem is to determine whether there exists a set of clusters, say  $\mathcal{X}$ , so that  $\bigcup_{C \in \mathcal{X}} C = \mathcal{N}$ , and  $\sum_{C \in \mathcal{X}} f(C) \geq \lambda$  where  $\lambda$  is a real number. We have the following theorem on the problem's complexity.

**Theorem 2.1.** *The robust clustering problem in CRN is NP-hard, when  $\delta_1 = 2$  and  $\delta_2 > 3$ .*

The proof is in Appendix ??.

### 3. CENTRALIZED SOLUTION FOR ROBUST CLUSTERING

When the global knowledge of the CRN is available to us, we can propose a centralized scheme as comparison. We obtain the set of  $\mathcal{S}$  which contains all the clusters in  $\mathcal{N}$ , i.e.,  $\mathcal{S} = \{C_1, C_2, \dots, C_i, \dots, C_{|\mathcal{S}|}\}$ <sup>§</sup> and there

<sup>§</sup>The subscript  $i$  means the  $i$ -th cluster in  $\mathcal{S}$ .

is  $\bigcup_{1 \leq i \leq |\mathcal{S}|} C_i = \mathcal{N}$ . The proposed centralized solution formulates the problem in Definition ?? as an optimization problem which is solved with standard software packages. The optimization decides on the clusters according to the following optimization formulation,

$$\begin{aligned} \max_{y_i, x_{ij}} \quad & \sum_{j=1}^N \sum_{i=1}^M (y_i \cdot t_{ij}) \\ \text{subject to} \quad & \sum_{i=1}^M x_{ij} = 1, \text{ for } \forall j = 1, \dots, N \\ & \sum_{j=1}^N x_{ij} = |C_i| \cdot y_i, \text{ for } \forall i = 1, \dots, M \\ & i \in \{1, 2, \dots, M\}, \quad j \in \{1, 2, \dots, N\} \end{aligned} \quad (1)$$

This problem is a binary linear programming problem, which can be solved by many available solvers.  $y_i$  and  $x_{ij}$  are two binary variables. Being either 1 or 0,  $y_i$  denotes whether the  $i$ -th cluster  $C_i$  in  $\mathcal{S}$  is chosen or not.  $x_{ij}$  indicates whether the CR node  $j$  resides in the cluster  $C_i$ , i.e.,  $x_{ij} = 1$  means node  $j$  resides in the cluster  $C_i$ .  $N$  is the total number of CR users in network  $\mathcal{N}$ ,  $M$  is the number of clusters in  $\mathcal{S}$ .

The constraints guarantee to obtain the clusters which together include all the CR users and don't overlap. The first constraint regulates that a CR node should reside in exactly one cluster. The second constraint regulates that when the  $i$ -th cluster  $C_i$  is chosen, there will be exactly  $|C_i|$  CR nodes residing in  $C_i$ .

The objective is to maximize the sum of the numbers of CCs in the clusters which constitute the CRN.  $t_{ij}$  is a constant and there is

$$t_{ij} = \frac{q_{ij}}{|C_i|} - p_i(C_i) \quad (2)$$

where constant  $q_{ij} = |K(C_i)|$  when node  $j \in C_i$ , and  $q_{ij} = 0$  when node  $j \notin C_i$ .  $p_i(C_i)$  is the size-related weight, which reflects the deviation of  $C_i$ 's size from the desired size. Assuming  $\delta$  is the desired size, then the weight  $p$  is decided according to the different cluster sizes, i.e.,  $1, 2, \dots, \sigma$ .

$$p(C_i) = \begin{cases} 0 & \text{if } |C_i| = \delta \\ \rho_1 & \text{if } ||C_i| - 1| = \delta \\ \rho_2 & \text{if } ||C_i| - 2| = \delta \\ \vdots & \\ \rho_\sigma & \text{if } ||C_i| - \sigma| = \delta \end{cases}$$

where  $\rho_1, \rho_2, \dots, \rho_\sigma$  are positive values and these is  $\sigma > \rho_2 > \rho_1 > 0$ .

When  $t_{ij}$  is replaced by  $\frac{q_{ij}}{|C_i|} - p(C_i)$ , the objective function becomes,

$$\max_{y_i, x_{ij}} \quad \sum_{j=1}^N \sum_{i=1}^M (y_i \cdot \frac{q_{ij}}{|C_i|} - y_i \cdot p(C_i))$$

The sum of the first items is the sum of CCs of all the chosen clusters. As to the second item, when  $w_i$  is 1 ( $C_i$  is chosen) and  $|C_i| \neq \delta$ , it will be negative, which

contradicts the direction of the optimization. Thus the second item discourages the appearance the clusters whose sizes deviate from  $\delta$ .

The difficulty of using this method lies in obtaining the set  $\mathcal{S}$ . In the worst case, i.e., the CRN forms a full connected graph, the size of  $\mathcal{S}$  is  $\sum_{r=1}^N \binom{N}{r} = 2^N - 1$ . Another obstacle comes from the fact that the centralized controller needs to be reliable at anytime, which is a challenge for CRN as the spectrum on the controller can not be guaranteed.

#### 4. DISTRIBUTED CLUSTERING ALGORITHM: ROSS

In this section we introduce the distributed clustering scheme ROSS. With ROSS, CR nodes form clusters based on the proximity of the available spectrum in their neighborhood after a series of interactions with their neighbors. ROSS consists of two cascaded phases: *cluster formation* and *membership clarification*. In the first phase, clusters are formed quickly and every CR user becomes either a cluster head or a cluster member, besides, cluster size control is implemented in this phase. In the second phase, non-overlapping clusters are formed in a way that the CCs of relevant clusters are mostly increased.

##### 4.1. Phase I - Cluster Formation

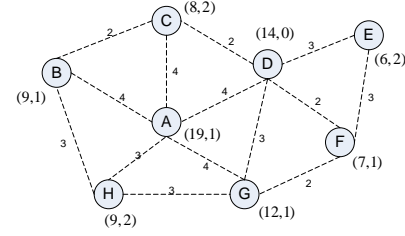
We assume that before conducting clustering, spectrum sensing, neighbor discovery and exchange of spectrum availability have been completed, so that every CR node is aware of the available channels for themselves and their neighbors. In this phase, cluster heads are determined after a series of comparisons with their neighbors. Two metrics are proposed to characterize the proximity in terms of available spectrum between CR node  $i$  and its neighborhood, which are used in the comparisons to decide on the cluster heads.

- *Individual connectivity degree*  $d_i$ :  $d_i = \sum_{j \in \text{Nb}(i)} |K_i \cap K_j|$ .  $d_i$  is the total number of the CCs between node  $i$  and each of its neighbors.
- *Neighborhood connectivity degree*  $g_i$ : the number of CCs which are available for  $i$  and all of its neighbors.  $g_i = |\bigcap_{j \in \text{Nb}(i) \cup i} K_j|$ , which represents the ability of  $i$  to form a robust cluster with its neighbors.

Individual connectivity degree  $d_i$  and neighborhood connectivity degree  $g_i$  together form the *connectivity vector*. Figure ?? illustrates an example CRN where every node's connectivity vector is shown. In Figure ??, some primary users are randomly deployed in the network, and they are free to choose one channel out of the set of licensed channels  $\mathcal{K}$  to operate. As to the secondary users, the channels they are allowed to use are those which are not ruled out by the primary users in their vicinities. The sets of the

indices of the available channels sensed by each node are:

$K_A = \{1, 2, 3, 4, 5, 6, 10\}$ ,  $K_B = \{1, 2, 3, 5, 7\}$ ,  $K_C = \{1, 3, 4, 10\}$ ,  $K_D =$



**Figure 1.** Connectivity graph of the example CRN and the connectivity vector  $(d_i, g_i)$  for each node. Primary users are not shown. The desired cluster size  $\delta = 3$ . The sets of the indices of the available channels sensed by each node are:  $K_A = \{1, 2, 3, 4, 5, 6, 10\}$ ,  $K_B = \{1, 2, 3, 5, 7\}$ ,  $K_C = \{1, 3, 4, 10\}$ ,  $K_D = \{1, 2, 3,$  Dashed edge indicates the end nodes are within each other's transmission range.

##### 4.1.1. Determining Cluster Heads and Forming Clusters

The procedure of determining the cluster heads is as follows. Each CR node decides whether it is a cluster head by comparing its connectivity vector with its neighbors. When CR node  $i$  has lower individual connectivity degree than all of its neighbors except for those which have already identified to be cluster heads, node  $i$  becomes a cluster head. If there is another CR node  $j$  in its neighborhood, which has the same individual connectivity degree as  $i$ , i.e.,  $d_j = d_i$  and  $d_j < d_k, \forall k \in \text{Nb}(j) \setminus \{\Lambda \cup i\}$  where  $\Lambda$  denotes the cluster heads, then the node between  $i$  and  $j$ , which has higher neighborhood connectivity degree will become the cluster head. If  $g_i = g_j$  as well, the node ID is used to break the tie, i.e., the one with smaller node ID becomes the cluster head. The node which is identified as a cluster head broadcasts a message to notify its neighbors of this change, and its neighbors which are not cluster heads become cluster members<sup>¶</sup>. The pseudo code for the cluster head decision and the initial cluster formation is shown in Algorithm ?? in the appendix.

After receiving the notification from a cluster head, a CR node  $i$  is aware that it becomes a member of a cluster. Consequently,  $i$  sets its individual connectivity degree to a positive number  $M > |\mathcal{K}| \cdot N$ , and broadcasts the new individual connectivity degree to all of its neighbors. When a CR node  $i$  is associated to multiple clusters, i.e.,  $i$  has received multiple notifications from different cluster heads,  $d_i$  is still set to be  $M$ . The manipulation of the individual connectivity degree of the cluster members accelerates the decision on the cluster heads. We have the following theorem to show that every secondary user will eventually

<sup>¶</sup> The reason for the occurrence of the cluster heads in the neighborhood of a new cluster head will be explained in Section ?? and ??)

be either integrated into a certain cluster or becomes a cluster head.

**Theorem 4.1.** *Given a CRN, it takes at most  $N$  steps ~~that before every secondary user either becomes cluster head, or gets included into at least one cluster.~~*

Here, by *step* we mean ~~one secondary user executing~~ Algorithm ?? ~~is executed by a secondary user~~ for one time. ~~The Proof is~~ Note that the minimum number of steps is 1, e.g., when a cluster head is elected and all the other nodes are within its neighborhood. The Proof on the upper bound of the steps is given in Appendix ?? ~~The procedure of the proof also illustrates the maximal time needed to conduct Algorithm ??.~~ Consider an extreme scenario, where all the secondary nodes sequentially execute Algorithm ??, i.e., they constitute a ~~list-chain~~ as discussed in the example in the proof. Assuming there is a token in the CRN, one node is allowed to execute Algorithm ?? only when it holds the token, and it then transfers the token to the neighboring node on the chain. If one step can be finished within a certain time span  $T$ , then the total time needed for the network to conduct Algorithm ?? is  $N * T$ . As Algorithm ?? can be executed concurrently by different secondary users, the needed time can be considerably reduced. If we apply Algorithm ?? to the example shown in Figure ??, ~~then~~ the outcome is shown in Figure ?? ~~Node B and H have the same individual connectivity degree, i.e.,  $d_B = d_H$ . As  $g_H = 2 > g_B = 1$ , node H becomes the cluster head and cluster  $C(H)$  is  $\{H, B, A, G\}$ .~~

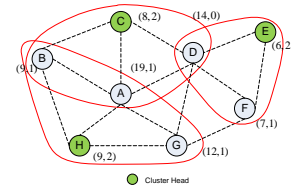
#### 4.1.2. The Existence of Common Channels

After executing Algorithm ??, certain formed clusters may not possess any CCs. As decreasing cluster size increases the CCs within a cluster, for those clusters having no CCs, certain nodes need to be eliminated to obtain at least one CC. The sequence of elimination is performed according to an ascending list of nodes which are sorted by the number of common channels between the nodes and the cluster head. In other words, the cluster member which has the least common channels with the cluster head is excluded first. If there are multiple nodes having the same number of common channels with the cluster head, the node whose elimination brings in more common channels will be excluded. If this criterion meets a tie, the tie will be broken by deleting the node with smaller node ID. It is possible that the cluster head excludes all of its neighbors, resulting in a singleton cluster which is composed by itself. The pseudo code for this procedure is shown in Algorithm ?? ~~As to the nodes which are eliminated from the previous clusters, they restore their original individual connectivity degrees, then execute Algorithm ?? and become either cluster heads or get included into other clusters afterwards according to Theorem ??.~~

During Phase I, whenever a CR node is decided to be a cluster head and accordingly forms a cluster, or its cluster's composition is changed, the cluster head will broadcast the updated information about its cluster, which includes the sets of available channels on all its cluster members.

#### 4.1.3. Cluster Size Control in Dense CRN

It is necessary to control the cluster size when CRN becomes denser. Both analysis and simulation [?] show that when applying ROSS, after the clusters are saturated with the increase of network density, the cluster size increases linearly with the network density, thus certain measures are needed to curb this problem. This task falls upon the cluster heads. To control the cluster size, cluster heads prune their cluster members to reach the desired cluster size. The desired size  $\delta$  is decided based on the capability of the CR users and the tasks to be conveyed. As there are overlaps between neighboring clusters, the sizes of the clusters formed in this phase are larger than that of the finally formed clusters. Hence, a cluster head excludes some cluster members when the cluster size exceeds a certain threshold. We set the threshold as  $t \cdot \delta$ , where  $\delta$  is the desired cluster size. The constant parameter  $t$  is usually larger than 1 and is dependent on the network density and CR nodes' transmission range ~~and  $t > 1$ .~~ In particular,  $t$  is between 1 and the ratio between the average number of neighbors and the desired size. A smaller or bigger  $t$  may lead to more clusters with smaller or larger size than  $\delta$ . Because of  $t$ , the threshold is larger than the desired size, then there will be some nodes choosing to affiliate with other clusters in the following phases. In particular, the cluster head removes the cluster members sequentially according to the following principle, the absence of one cluster member leads to the maximum increase of the CCs within the cluster. This process ends when each cluster's size is smaller or equal to  $t \cdot \delta$ . This procedure is similar with that in Section ??, thus Algorithm ?? can be reused. ~~The  $t$  is set to 1.3.~~



**Figure 2.** Clusters formation after the phase I of ROSS. Nodes A, B, D are debatable nodes as they belong to multiple clusters.

#### 4.2. Phase II - Membership Clarification

As to the example CRN shown in Figure ??, the resulted clusters are shown in Figure ?? after running phase I of ROSS. We notice that nodes A, B, D are included in more than one cluster. We refer to these nodes as *debatable nodes* as their cluster affiliations are not decided. The clusters which include the debatable node  $i$  are called *claiming clusters* of node  $i$ , and the set of these clusters is denoted as  $S_i$ . The debatable nodes which are generated from the first phase of ROSS should be exclusively associated with only one cluster and be removed from the other claiming clusters, this procedure is called *cluster membership clarification*.

#### 4.2.1. Distributed Greedy Algorithm (DGA)

Assuming a debatable node  $i$  which needs to decide one cluster  $C \in S_i$  to stay and leaves the other clusters in  $S_i$ , then the principle for  $i$  is its decision should result in the greatest increase of CCs in all its claiming clusters. As node  $i$  has been notified of the spectrum availability on all the nodes in each claiming cluster, node  $i$  is able to calculate how many more CCs will be produced in a claiming cluster if  $i$  leaves that cluster. Then node  $i$  decides on the cluster  $C \in S_i$ , if  $i$  leaving cluster  $C$  results in less increased CCs than leaving any other claiming clusters in  $S_i$ . When there comes a tie between two claiming clusters,  $i$  chooses to stay in the cluster whose cluster head shares the most CCs with  $i$ . When the tie still exists, node  $i$  chooses to stay in the claiming cluster which has the smallest size. Node IDs of cluster heads will be used to break tie if all the previous metrics could not decide on the unique claiming cluster for  $i$  to stay. The pseudo code of this algorithm is given in Algorithm ???. After deciding its membership, debatable node  $i$  notifies all its claiming clusters of its choice, and the claiming clusters from which node  $i$  leaves also broadcast their new cluster composition and the spectrum availability on all their cluster members.

The autonomous decisions made by the debatable CR nodes raise the concern on the endless chain effect in the membership clarification phase. A debatable node's choice is dependent on the compositions of its claiming clusters, which can be changed by other debatable nodes' decisions. As a result, the debatable node which makes decision first may change its original choice, and this process may go on forever. To erase this concern, we formulate the process of membership clarification into a game, where an equilibrium is reached after a finite number of best response updates made by the debatable nodes.

#### 4.2.2. Bridging ROSS-DGA with Congestion Game

Game theory is a powerful mathematical tool for studying, modeling and analyzing the interactions among individuals. A game consists of three elements: a set of players, a selfish utility for each player, and a feasible strategy space for each player. In a game, the players are rational and intelligent decision makers, which are related with one explicit formalized incentive expression (the utility or cost). Game theory provides standard procedures to study its equilibriums [?]. In the past few years, game theory has been extensively applied to problems in communication and networking [?, ?]. Congestion game is an attractive game model which describes the problem where participants compete for limited resources in a non-cooperative manner, it has the good property that Nash equilibrium can be achieved after finite steps of best response dynamic, i.e., each player chooses the strategy to maximize/minimize its utility/cost with respect to the other players' strategies. The framework of the congestion

game has been used to model certain problems in internet-centric applications or cloud computing, where self-interested clients compete for the centralized resources and meanwhile interact with each other. For example, server selection is involved in distributed computing platforms [?], or users downloading files from cloud, etc.

To formulate the debatable nodes' membership clarification into the desired congestion game, we reexamine this process from a different/opposite perspective. From the new perspective, the debatable nodes are not included in any cluster and they need to decide on one cluster to join. When a debatable node  $i$  join one cluster  $C$ , the decrease of CCs in cluster  $C$  is  $\sum_{C \in S_i} \Delta|K(C)| = \sum_{C \in S_i} (|K(C)| - |K(C \cup i)|)$ . Then, node  $i$  chooses the cluster  $C$ , where the decrease of CCs in cluster  $C$  is smaller than the decrease if  $i$  would have joined any other claiming cluster in  $S_i$ . The relation between the debatable nodes and the claiming clusters is shown in Figure ??.

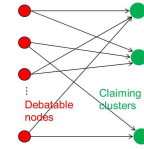


Figure 3. Debatable nodes and claiming clusters

In the following, we show that the decision of debatable nodes to clarify their membership can be mapped to the behaviour of the players in a *player-specific singleton congestion game* when proper cost function is given. The game to be constructed is represented with a 4-tuple  $\Gamma = (\mathcal{P}, \mathcal{R}, \sum_{i \in \mathcal{P}}, f)$  with the following elements:

- $\mathcal{P}$ , the set of players in the game, which are the debatable nodes in our problem.
- $\mathcal{R} = \cup S_i, i \in \mathcal{P}$ , the set of the resources for players to choose. In our problem,  $S_i$  is the set of the claiming clusters of  $i$ , and  $\mathcal{R}$  is the set of all claiming clusters.
- Strategy space  $\sum_i, i \in \mathcal{P}$ ,  $\sum_i$  is the set of the claiming clusters  $S_i$ . As debatable node  $i$  is supposed to choose only one claiming cluster, then only one piece of resource will be allocated to  $i$ .
- The utility (cost) function  $f(C)$  as to a resource  $C$ .  $f(C) = \Delta|K^i(C)|, C \in S_i$ , which represents the decreased number of CCs in cluster  $C$  when debatable node  $i$  joins  $C$ . As to cluster  $C \in S_i$ , the decrease of CCs caused by including the debatable nodes is  $\sum_{i: C \in S_i, i \rightarrow C} \Delta|K^i(C)|$ .  $i \rightarrow C$  means  $i$  joins cluster  $C$ . Obviously this function is non-decreasing with respect to the number of nodes joining cluster  $C$ .

The utility function  $f$  is not purely decided by the number of players accessing the resource (debatable nodes join claiming clusters), which happens in a canonical congestion game. The reason is in this game the channel availability on debatable nodes is different.



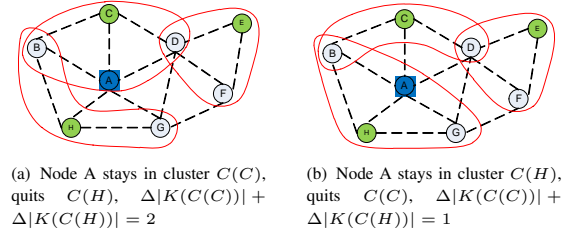
Given two same groups of debatable nodes and their sizes are the same, when the nodes are not completely the same (neither are the channel availabilities on these nodes), the cost happened on one claiming cluster could be different if the two groups of debatable nodes join that cluster respectively. Hence, this congestion game is player specific [?]. In this game, every player greedily updates its strategy (choosing one claiming cluster to join) if joining a different claiming cluster minimizes the decrease of CCs  $\sum_{i: C \in S_i} \Delta |K^i(C)|$ , and a player's strategy in the game is exactly the same with the behaviour of a debatable node in the membership clarification phase.

As to singleton congestion game, there exists a pure equilibria which can be reached with the best response update, and the upper bound for the number of steps before convergence is  $n^2 * m$  [?], where  $n$  is the number of players, and  $m$  is the number of resources. In our problem, the players are the debatable nodes, and the resources are the claiming clusters. Thus the number of steps can be expressed as  $\mathcal{O}(N^3)$ . In fact, the upper bound for the number of steps which are involved in this process is much smaller than  $N^3$ . The percentage of debatable nodes in the network is shown in Figure ??, which is between 10% to 60% of the total CR nodes in the network. On the other hand, the number of clusters heads is dependent on the network density and the CR node's transmission range as mentioned in Section ?. The simulation in [?] shows the cluster heads are only 3.4% to 20% of the total CR nodes.

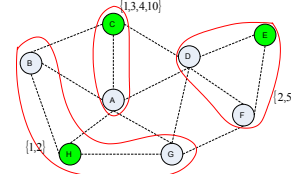
#### 4.2.3. Distributed Fast Algorithm (DFA)

On the basis of ROSS-DGA, we propose a faster version ROSS-DFA which differs from ROSS-DGA in the second phase. With ROSS-DFA, debatable nodes decide their respective cluster heads only once. The debatable nodes consider their claiming clusters to include all their debatable nodes, thus the membership of claiming clusters is static and all the debatable nodes can make decisions simultaneously without considering the change of membership of their claiming clusters. As ROSS-DFA is quicker than ROSS-DGA, the former is especially suitable for the CRN where the channel availability changes frequently. To run ROSS-DFA, debatable nodes execute only one loop in Algorithm ?.

Now we apply both ROSS-DGA and ROSS-DFA to the ~~toy~~ network in Figure ?? which has been applied the phase I of ROSS. In the network, node A's claiming clusters are cluster  $C(C)$ ,  $C(H) \in S_A$ , their members are  $\{A, B, C, D\}$  and  $\{A, B, H, G\}$  respectively. The two possible strategies of node A is illustrated in Figure ?. In Figure ??, node A staying in  $C(C)$  and leaving  $C(H)$  brings 2 more CCs to  $S_A$ , which is more than that brought by another strategy shown in ?. After the decisions made similarly by the other debatable nodes B and D, the final clusters are formed as shown in Figure ?.



**Figure 4.** Membership clarification: possible cluster formations caused by node A's different choices



**Figure 5.** Final formation of clusters. Common channels are shown beside corresponding clusters.

## 5. PERFORMANCE EVALUATION

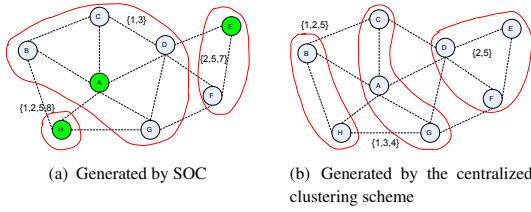
The schemes involved in the simulation are as follows,

- ROSS without size control: ROSS-DGA, ROSS-DFA.
- ROSS with size control, i.e., ROSS- $\delta$ -DGA and ROSS- $\delta$ -DFA where  $\delta$  is the desired cluster size. In the following, we refer to the above mentioned four schemes as the variants of ROSS.
- SOC [?], a distributed clustering scheme pursuing cluster robustness.
- Centralized robust clustering scheme. As shown in Section ??, the centralized robust clustering scheme is formulated as an integer linear optimization problem and is solved by MATLAB with the function *bintprog*.

The ROSS without size control mechanism is similar with the schemes proposed in [?]. The authors of [?] compared SOC with other schemes in terms of the average number of CCs of the formed cluster, on which SOC outperforms other schemes by 50%-100%. SOC's comparison schemes are designed either for ad hoc network without consideration of channel availability [?], or for CRN but just considering connection among CR nodes [?]. Thus SOC is the only distributed scheme as comparison. As to the CRN shown in Figure ??, the resulting clusters by the centralized scheme and SOC are shown in Figure ?.

We investigate the schemes ~~with respect to four metrics~~ in the following aspects.

- **The average number of CCs per non-singleton cluster.** Non-singleton cluster refers to the cluster whose cluster size is larger than one. Previous work [?] and [?]



**Figure 6.** Final clusters formed by the centralized clustering scheme and SOC.

claim that the larger average number of CCs over all the clusters indicates robustness, from which we see two flaws. First, the unclustered CR nodes (synonym of singleton clusters) should not be considered when calculating the average number of CCs, as singleton clusters don't contribute to the collaborative computing or sensing. Second, the average number of CCs doesn't necessarily indicate the robustness of individual clusters, because the ability for a cluster to sustain also depends on cluster size and the locations of the cluster members, but these information can not be illustrated in the average number of CCs. In the performance evaluation, we will examine the metric of average number of CCs per non-singleton cluster, which excludes the bias brought in by the unclustered CR nodes. Moreover, we will examine whether this metric reflects the robustness of the clusters.

- **Cluster sizes.** We investigate the distribution of CRs residing in the formed clusters with different sizes.
- **Robustness of the clusters against newly added PUs.** We increase Robustness is illustrated by the number of PUs to challenge the non-singleton clusters, and count the number of the unclustered CR nodes the unclustered CR nodes in the CRN, after the CRN being challenged by the increasing number of PUs. This metric indicates the robustness of the clusters, i.e., as to the clusters formed for a given CRN and spectrum availability, how many CR nodes can still be benefited from the clusters when the spectrum availability decreases.
- **Cluster sizes.** We investigate the distribution of CRs residing in the formed clusters with different sizes.
- **Amount of control messages involved.** We investigate the number of control messages involved in the clustering process.
- **Influence from inaccurate spectrum sensing.** The above simulations are conducted with the assumption of perfect spectrum sensing. As the errors in spectrum sensing is inevitable, we are interested to see the performance of the distributed schemes when the spectrum sensing is not perfect. The false negative in spectrum sensing, which misdetects the presence of the active primary users, is harmful to the primary users

and should be avoided as much as possible. On the contrary, false positive i.e., which reports the presence of active primary users when there are actually not, only decreases the available spectrum for the secondary users. In this regard, we assume only the false negative exists in the spectrum sensing.

We assume when a secondary user is within the transmission range of an active primary user, the probability that it misdetects and thus regards a channel is available equals to the rate of false negative. The secondary users make clustering decisions based on the imperfect spectrum sensing. After the clustering process is completed, we correct the spectrum availability with the ground truth. Then certain formed clusters may be affected as their CCs which are obtained due to false negative will be revoked.

**Simulation**—The simulation consists of two parts, first we investigate the performance of centralized scheme and the distributed schemes in a small network, as there is no polynomial time solution available to solve the centralized problem. In the second part, we investigate the performance of the proposed distributed schemes in the CRN with different scales and densities. The following simulation setting is the same for both simulation parts. CRs and PUs are deployed on a two-dimensional Euclidean plane. The number of licensed channels is 10, each PU is operating on each channel with probability of 50%. The other parameters i.e., the number of CR and PU, and their transmission ranges are given in the beginning of the respective simulation sections. The constant  $t$  which is used to control cluster size for ROSS (discussed in Section ??) is 1.3. CR users are assumed to be able to sense the existence of primary users and identify available channels. All primary and CR users are assumed to be static during the process of clustering. The simulation is written in C++, and the performance results are averaged over 50 randomly generated topologies, and the confidence interval corresponds to 95% confidence level.

## 5.1. Centralized Schemes vs. Decentralized Schemes

There In this part of simulation, there are 10 primary users and 20 CR users dropped randomly (with uniform distribution) within a square area of size  $A^2$ , where  $A$  is a positive value and we set the transmission ranges of primary and CR users to  $A/3$ . When clustering scheme is executed, around 7 channels are available on each CR node. The desired cluster size  $\delta$  is 3. As for the centralized scheme, the parameters used in the punishment for choosing the clusters with undesired sizes are set as follows,  $\rho_1 = 0.4$ ,  $\rho_2 = 0.6$ .

### 5.1.1. CCs in Non-singleton Clusters

From Figure ??, we can see the centralized schemes outperform the distributed schemes. Among the distributed schemes, SOC achieves the most CCs. The reason is, SOC

is liable to group the neighboring CRs which share the most abundant spectrum together, no matter how many of them are there, thus the number of CC of the formed clusters is higher. On the other hand, SOC generates the most unclustered CRs. As to the variants of ROSS, we notice that the greedy mechanism increases CCs in non-singleton clusters significantly.

### 5.1.2. Cluster Size

Figure ?? depicts the empirical cumulative distribution of the CRs in clusters of different sizes, from which we have two conclusions. First, given the channel availability in the CRN, SOC generates more unclustered CR nodes than other schemes. The centralized schemes don't produce unclustered CR nodes in the simulation, the unclustered nodes generated by ROSS-DGA/DFA account for 3% of the total CR nodes, as comparison, 10% of nodes are unclustered when applying SOC. ROSS-DGA and ROSS-DFA with size control feature generate 5%-8% unclustered CR nodes, which is due to the cluster pruning procedure (discussed in section ?? and section ??). Second, the centralized schemes and cluster size control mechanism of ROSS generate clusters with the desired cluster size. As to ROSS-DFG and ROSS-DFA with size control feature, CR nodes reside averagely in clusters whose sizes are 2, 3 and 4. The sizes of clusters resulted from ROSS-DGA and ROSS-DFA are disperse, but appear to be better than SOC, i.e., the 50% percentiles for ROSS-DGA, ROSS-DFA and SOC are 4.5, 5, and 5.5, and the 90% percentiles for the three schemes are 8, 8, and 9, the corresponding sizes of ROSS are closer to the desired size.

### 5.1.3. Robustness of the formed clusters

In this part of simulation, we put PUs sequentially into CRN to decrease the available spectrum. 10 PUs are in the network in the beginning, then extra 19 batches of PUs are added sequentially, where each batch includes 5 PUs. Figure ?? shows certain clusters can not maintain and the number of unclustered CR nodes grows when the number of PUs increases. The centralized scheme with desired size of 2 generates the most robust clusters, meanwhile, SOC results in the most vulnerable clusters. The centralized scheme with desired size of 3 doesn't outperform the variants of ROSS, because pursuing cluster size prevents forming the clusters with more CCs. In contrary, the variants of ROSS generate some smaller clusters which are more likely to maintain when there are more PUs.

The above observation shows that the average number of CCs of non-singleton clusters doesn't necessarily illustrate the robustness of cluster, i.e., SOC obtains the most CCs for the clusters which are meanwhile the most vulnerable. Besides, with similar distribution of sizes, the clusters generated by ROSS-DGA and ROSS-DFA are more robust than that by SOC.

### 5.1.4. Control Signaling Overhead

In this section we compare the overhead of signaling involved in different clustering schemes. We count the number of *transmissions of control messages* as message complexity [?], and without distinguishing broadcast or uni-cast control messages. In Section ??, this metric is synonymous with the *the number of updates*.

As to ROSS, in the first phase the maximal number of broadcast is  $N$  according to ?. The upper bounds for the transmissions are  $n^2m$  and  $n$  for ROSS-DGA and ROSS-DFA respectively. Scheme SOC consists of three rounds, and in each round every node needs to broadcast to do comparisons and cluster mergers. The centralized scheme is conducted at the centralized control device, which involves information aggregation and clustering decision dissemination. We adopt the backbone structure proposed in [?] to analyze the centralized scheme's message complexity. We apply ROSS to generate cluster heads which serve as the backbone. In the process of information aggregation, all the nodes transmit information to the cluster heads which forward the messages to the controller, then in the process of dissemination, all the cluster heads and the debatable nodes broadcast the clustering result, thus the upper bound for the number of broadcast is  $N + m + n$ .

The number of control messages which are involved in ROSS variants and the centralized scheme is related with the number of debatable nodes. Figure ?? shows the percentage of debatable nodes with different network densities. Table ?? shows the message complexity, quantitative amount of the control messages, and the size of control messages. Figure ?? shows the analytical result of the amount of transmissions involved in different schemes.

## 5.2. Comparison among the Distributed Schemes

In this section we investigate the performances of the proposed distributed clustering schemes in CRN with different network scales and densities. The transmission range of CR is  $A/5$ , PU's transmission range is  $2A/5$ . The initial number of PU is 30. The desired sizes adopted are listed in the Table ??, which is about 60% of the average number of neighbors.

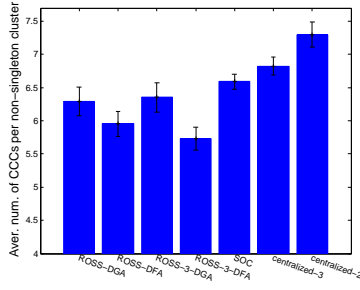
### 5.2.1. Number of CCs per Non-singleton Clusters

The average number of CCs of the non-singleton clusters is shown in Figure ?. SOC achieves the most CCs per non-singleton cluster, but the lead over the variants of ROSS decreases significantly when  $N$  increases.

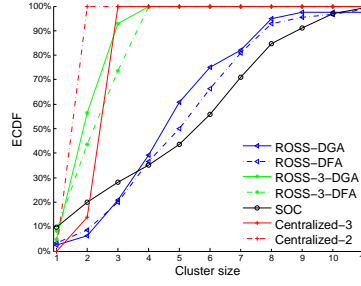
### 5.2.2. Robustness of the formed clusters

Here we see how vulnerable the clusters are when they are exposed to the increasing influence of the PUs. We increase the primary users' activity by importing 20

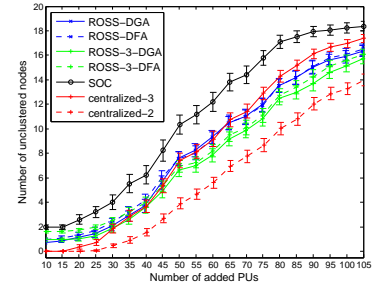




**Figure 7.** Average number of CCs of non-singleton clusters



**Figure 8.** Cumulative distribution of CRs residing in clusters with different sizes



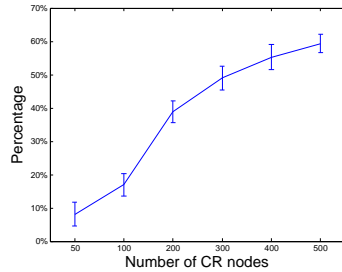
**Figure 9.** Number of unclustered CRs with decreasing spectrum availability

**Figure 10.** Comparison between the distributed and centralized clustering schemes ( $N = 20$ )

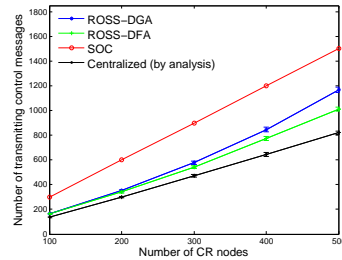
**Table II.** Signalling overhead

Scheme	Message Complexity	Quantitative number of messages	Content and size of the message
ROSS-DGA, ROSS- $\delta$ -DGA	$\mathcal{O}(N^3)$ (worst case)	$N + n^2m$ (upper bound)	PhaseI: ID, $d_i, g_i$ , which are 3 bytes; PhaseII: Cluster head $i$ broadcasts channel availability to all members, where are $ C(i)  \mathcal{K} $ bytes
ROSS-DFA, ROSS- $\delta$ -DFA	$\mathcal{O}(N)$ (worst case)	$N + n$ (upper bound)	
SOC	$\mathcal{O}(N)$	$3N$	Every CR node $i$ broadcasts channel availability on all cluster members, which is $ C(i)  \mathcal{K} $ bytes
Centralized	$\mathcal{O}(N)$	$N + n + m$ (upper bound)	clustering result, which is $2N$ bytes ??

<sup>a</sup> Assuming the data structure of the clustering result is in the form of  $\{i, C\}$ ,  $i \in C$ ,  $i \in \mathcal{N}$ .



**Figure 11.** The percentage of debatable nodes after phase I of ROSS.

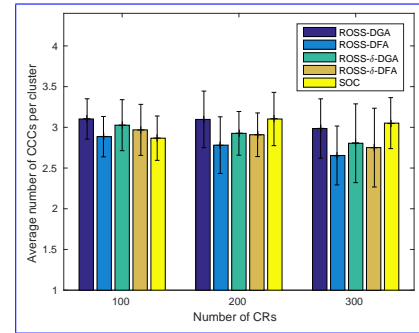


**Figure 12.** Quantitative amount of control messages.

batches of PUs sequentially in the CRN, where each batch includes 10 PUs. Figure ?? and ?? show that when  $N = 100$  and 200, compared with the variants of ROSS, more

**Table III**

Number of CRs	100	200	300
Average num. of neighbors	9.5	20	31
Desired size $\delta$	6	12	20



**Figure 13.** Average number of CCs of non-singleton clusters

unclustered CR nodes are ~~generated by SOC~~ caused by SOC with the same intensity of PUs' activities. When  $N = 300$  as shown in Figure ?? and the new PUs are not many, ROSS-DGA/DFA generate slightly more unclustered CR nodes than SOC, but SOC's performance deteriorates quickly when the PUs continue increasing. From Figure ?? to ??, we can see that significantly less unclustered CR

nodes are generated by the variants of ROSS which have size control mechanism. Besides, the greedy mechanism moderately strengthens the robustness of the clusters. We only show the average values of the variants of ROSS as their confidence intervals overlap.

### 5.2.3. Cluster Size Control

Figure ?? shows when the network density increases, i.e.,  $N$  changes from 100 to 300, the number of generated clusters by SOC increases linearly, whereas that by ROSS increases by a smaller margin. This result coincides with the analysis in Section ?. To better understand the distribution of the sizes of formed clusters, for each network density, we depict the empirical cumulative distribution of CR nodes which are in clusters with different sizes in Figures ?? ?? ?? respectively.

The variants of ROSS generate more clusters whose sizes are closer to the desired size, i.e., when  $N = 100$  and desired cluster size is 6 as shown in Figures ??, 90% of CR nodes are in the clusters whose sizes are from 3 to 9, while as to SOC, only 17% of nodes are in the clusters with these sizes. Similarly, when  $N = 200$  and the desired size is 12 as shown in Figure ??, 80% of nodes are in the clusters whose sizes are from 6 to 18, meanwhile only 30% of nodes constitute clusters of these sizes when SOC is executed. The clusters sizes from ROSS- $\delta$ -DGA and ROSS- $\delta$ -DFA concentrates more than that from ROSS-DGA and ROSS-DFA. In contrary, the clusters from SOC demonstrates obvious divergence on cluster sizes.

The limitation of distributed scheme ROSS is it doesn't generate clusters whose sizes exceed the cluster head's neighborhood. The reason is with ROSS, cluster heads form clusters on the basis of their neighborhood, and don't involve the nodes which are outside the neighborhood.

### 5.2.4. The Performance with False Negative in Spectrum Sensing

Figure ?? shows the average number of CCs decrease slightly when the false negative rate increases. The size distribution of the ROSS-DGA, ROSS- $\delta$ -DGA and SOC is shown in Figure ?. For all the schemes, when false negative increases, the number of singleton clusters increases accordingly. The clusters from SOC are affected by the sensing errors greatly, i.e., when false negative rate is 30%, a lot more small clusters are generated than that when other false negative rates present. SOC demonstrates obvious divergence when certain rate of false negative exists, in contrary, ROSS variants are resilient to the false negative in terms of cluster sizes. The negotiation within the neighborhoods adopted by ROSS variants rules out the channel which is due to false negative successfully, which doesn't work well in SOC.

### 5.3. Insights Obtained from the Simulation

The simulation made-with-the-with large CRN network confirms that-made-with-the-the conclusion drawn from the small CRN network, which is that the-. First, different

from the assumption adopted in all the previous works, the average number of CCs along-doesn't-per cluster can tell the robustness of the clusters, because the cluster size and the constitution of the cluster also affect the robustness.

The-against the increasing influence of the primary users. Second, the centralized clustering scheme is able to form the clusters which satisfy the requirement on cluster size strictly, and the clusters are robust against the-PU's-PU's' increasing activity, besides, it involves the smallest control overhead in the process of clustering. As-Third, as distributed schemes, the variants of ROSS outperform SOC considerably on-three-metrics. First, the-in four aspects.

- The variants of ROSS generate less unclustered nodes than SOC for a given CRN, and the resulted clusters are more robust than SOC when PUs become more active. Second, the-
- The signaling overhead involved in ROSS is about half of that needed for SOC, and the signaling messages are much shorter than the latter. Third, the-
- The sizes of the clusters generated by ROSS demonstrate smaller discrepancy than that of SOC.
- The variants of ROSS are more resilient against the imperfect spectrum sensing.

Moreover, the ROSS variants with size control features achieve similar performance to the centralized scheme in terms of cluster size, and the cluster robustness is similar when applying the variants of ROSS and the centralized scheme respectively. As-to-Among the variants of ROSS, the greedy mechanism in ROSS-DGA helps to improve the performance on cluster size and cluster robustness at the cost of increased signaling overhead.

## 6. CONCLUSION

In this paper we investigate the robust clustering problem in CRN extensively and propose both centralized and distributed clustering solutions. We give the mathematical description of the problem and prove the NP hardness of it. The proposed clustering schemes generate clusters which have long life expectancy against the primary users' activity, and the generated clusters have similar sizes with the desired one. Through simulation, the distributed schemes demonstrate similar performance with the centralized scheme in terms of cluster robustness, signaling overhead and cluster sizes, and outperform the comparison distributed scheme on all metrics.

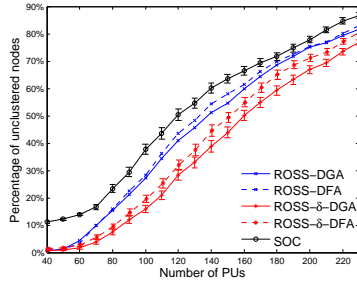


Figure 14. 100 CRs

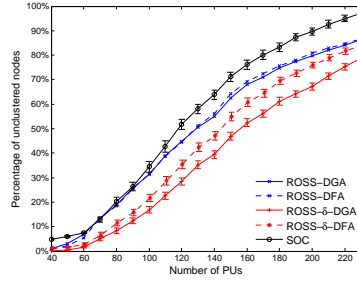


Figure 15. 200 CRs

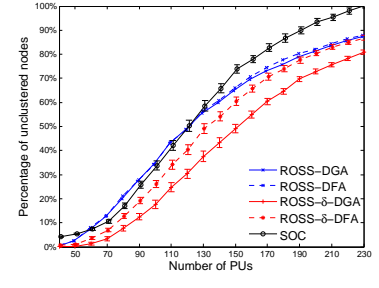


Figure 16. 300 CRs

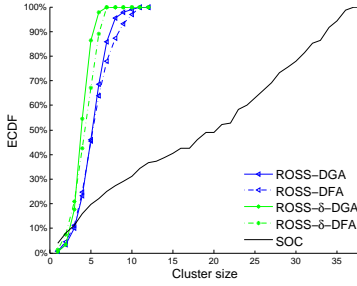


Figure 17. 100 CRs, 30 PUs in network

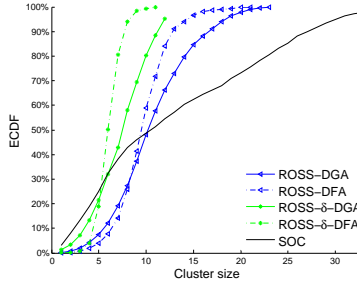


Figure 18. 200 CRs, 30 PUs in network

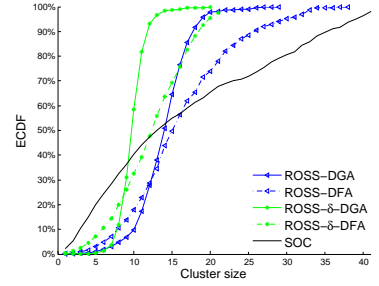


Figure 19. 300 CRs, 30 PUs in network

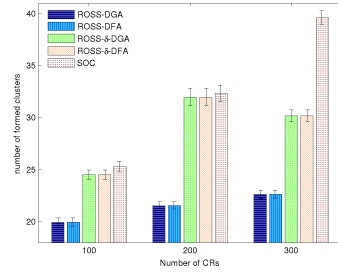


Figure 20. The number of formed clusters.

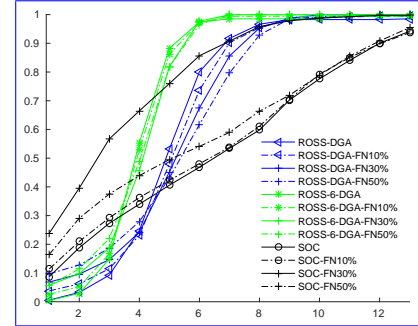


Figure 22. 100 CRs with false negative in spectrum sensing, 30 PUs in network

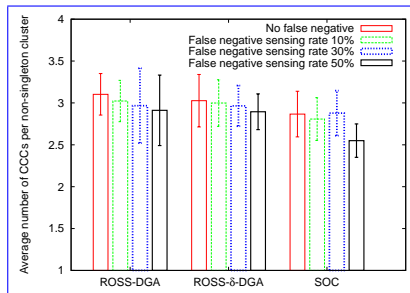


Figure 21. The number of CCs per non-singleton cluster with the presence of spectrum sensing false negative

## Appendices

### A. PEUDO CODE FOR THE ALGORITHM ??, ?? AND ??

Trans. Emerging Tel. Tech. 2017; 00:??-?? © 2017 John Wiley & Sons, Ltd.

B. PROOF OF THEOREM ??  
Prepared using ettauth.cls

Proof

We consider a CRN which can be represented as a

connected graph. To simplify the discussion, we assume the secondary users have unique individual connectivity degrees. Each user has an identical ID and a neighborhood connectivity degree. This assumption is fair as the neighborhood connectivity degrees and node ID are used to break ties in Algorithm ??, when the individual connectivity degrees are unique, it is not necessary to use the former two metrics.

For the sake of contradiction, let us assume there exist a secondary user  $\alpha$  which is not included into any cluster. Then there is at least one node  $\beta \in \text{Nb}(\alpha)$  such that  $d_\alpha > d_\beta$  (otherwise  $\alpha$  becomes cluster head). In this case, according to Algorithm ??,  $\beta$  is not included in any clusters, because otherwise  $d_\beta = M$ , a large

---

**Algorithm 1:** ROSS phase I: cluster head determination and initial cluster formation for CR node  $i$ 


---

**Input:**  $d_j, g_j, j \in \text{Nb}(i) \setminus \Lambda$ ,  $\Lambda$  denotes the set of cluster heads among  $\text{Nb}(i)$ . Empty sets  $\tau_1, \tau_2$

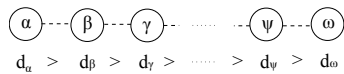
**Result:** Returning 1 means  $i$  is cluster head, and  $d_j$  is set to 0,  $j \in \text{Nb}(i) \setminus \Lambda$ . Returning 0 means  $i$  is not cluster head.

```

1 if  $\nexists j \in \text{Nb}(i) \setminus \Lambda$ , such that  $d_i \geq d_j$  then
2   return 1;
3 end
4 if  $\exists j \in \text{Nb}(i) \setminus \Lambda$ , such that  $d_i > d_j$  then
5   return 0;
6 else
7   if  $\nexists j \in \text{Nb}(i) \setminus \Lambda$ , such that  $d_j == d_i$  then
8      $\tau_1 \leftarrow j$ 
9   end
10 end
11 if  $\nexists j \in \tau_1$ , such that  $g_i \leq g_j$  then
12   return 1;
13 end
14 if  $\exists j \in \tau_1$ , such that  $g_i < g_j$  then
15   return 0;
16 else
17   if  $\nexists j \in \tau_1$ , such that  $g_j == g_i$  then
18      $\tau_2 \leftarrow j$ 
19   end
20 end
21 if  $ID_i$  is smaller than any  $ID_j, j \in \tau_2 \setminus i$  then
22   return 1;
23 end
24 return 0;
```

---

positive integer, which contradicts to  $d_\alpha > d_\beta$ . Now, we distinguish between two cases: If  $\beta$  becomes cluster head, node  $\alpha$  is included, the assumption is not true. If  $\beta$  is not a cluster head, then  $\beta$  is not in any cluster, we can repeat the previous analysis made on node  $\alpha$ , and deduce that node  $\beta$  has at least one neighboring node  $\gamma$  with  $d_\gamma < d_\beta$ . So far, when there is no cluster head identified, the unclustered nodes, i.e.,  $\alpha, \beta$  form a linked list, where their connectivity degrees monotonically decrease. But this list will not continue to grow, because the minimum individual connectivity degree is zero, and the length of this list is upper bounded by the total number of nodes in the CRN. An example of the formed node series is shown as Figure ??.



**Figure 23.** The node series discussed in the proof of Theorem ??, the deduction begins from node  $\alpha$

In this example, node  $\omega$  is at the tail of a list. As  $\omega$  does not have neighboring nodes with lower individual connectivity degree,  $\omega$  becomes a cluster head. Then  $\omega$

---

**Algorithm 2:** ROSS phase I: cluster head guarantees the availability of CC (start from line 1) / cluster size control (start from line 2)

---

**Input:** Cluster  $C$ , empty sets  $\tau_1, \tau_2$

**Output:** Cluster  $C$  has at least one CC, or satisfies the requirement on cluster size

```

1 while  $K_C = \emptyset$  do
2   while  $|C| > t \cdot \delta$  do
3     if  $\exists$  only one  $i \in C \setminus h(C)$ ,
4        $i = \arg \min(|K_{h(C)} \cap K_i|)$  then
5        $C = C \setminus i$ ;
6     else
7        $\exists$  multiple  $i$  which satisfies
8        $i = \arg \min(|K_{h(C)} \cap K_i|)$ ;
9        $\tau_1 \leftarrow i$ ;
10    end
11    if  $\exists$  only one  $i \in \tau_1$ ,
12       $i = \arg \max(|\cap_{j \in C \setminus i} K_j| - |\cap_{j \in C} K_j|)$ 
13      then
14         $C = C \setminus i$ ;
15      else
16         $C = C \setminus i$ , where  $i = \arg \min_{i \in \tau_1} ID_i$ 
17      end
18    end
19  end
20 end
```

---

incorporates all its one-hop neighbors (here we assume that every newly formed cluster has common channels), including the nodes which precede  $\omega$  in the list. The nodes which join a cluster set their individual connection degrees to  $M$ , which makes the node immediately precede in the list to become a cluster head. In this way, cluster heads are generated from the tail to the head in the list, and every node in the list is in at least one cluster, which contradicts the assumption that  $\alpha$  is not included in any cluster.

If we see a secondary user *becoming a cluster head*, or *becoming a cluster member* as one step, as the length of the list of secondary users is not larger than  $N$ , there are  $N$  steps for this scenario to form the initial clusters. □

## C. PROOF OF THEOREM ??

### Proof

To prove the robust clustering problem is NP-hard, we reduce the *maximum weighted  $k$ -set packing problem*, which is NP-hard when  $k \geq 3$  [?], to the the robust clustering problem to show the latter is at least as hard as the former. Given a collection of sets of cardinality at most  $k$  and with weights for each set, the maximum weighted packing problem is that of finding a collection of disjoint sets of maximum total weight. The decision version of the weighted  $k$ -set packing problem is,

---

**Algorithm 3:** Debatable node  $i$  decides its affiliation in phase II of ROSS

---

**Input:** all claiming clusters  $C \in S_i$   
**Output:** one cluster  $C \in S_i$ , node  $i$  notifies all its claiming clusters in  $S_i$  about its affiliation decision.

```

1 while  $i$  has not chosen the cluster, or  $i$  has joined
   cluster  $\tilde{C}$ , but  $\exists C' \in S_i, C' \neq \tilde{C}$ , which has
    $|K(C' \setminus i)| - |K(C')| < |K(C \setminus i)| - |K(C)|$  do
2   if  $\exists$  only one  $C \in S_i$ ,
      $C = \arg \min(|K(C \setminus i)| - |K(C)|)$  then
3     return  $C$ ;
4   else
5      $\exists$  multiple  $C \in S_i$  which satisfies
      $C = \arg \min(|K(C \setminus i)| - |K(C)|)$ ;
6      $\tau_1 \leftarrow C$ ;
7   end
8   if  $\exists$  only one  $C \in \tau_1$ ,
      $C = \arg \max(K_{h(C)} \cap K_i)$  then
9     return  $C$ ;
10  else
11     $\exists$  multiple  $C \in S_i$  which satisfies
     $C = \arg \max(K_{h(C)} \cap K_i)$ ;
12     $\tau_2 \leftarrow C$ ;
13  end
14  if  $\exists$  only one  $C \in \tau_2, C = \arg \min |C|$  then
15    return  $C$ ;
16  else
17    return  $\arg \min_{C \in \tau_2} h(C)$ ;
18  end
19 end

```

---

**Definition 2.** Given a finite set  $\mathcal{G}$  of non-negative integers where  $\mathcal{G} \subseteq \mathbb{N}$ , and a collection of sets  $\mathcal{Q} = \{S_1, S_2, \dots, S_m\}$  where  $S_i \subseteq \mathcal{G}$  and  $\max(|S_i|) \geq 3$  for  $1 \leq i \leq m$ . Every set  $S$  in  $\mathcal{Q}$  has a weight  $\omega(S) \in \mathbb{N}^+$ . The problem is to find a collection  $\mathcal{I} \subseteq \mathcal{Q}$  such that  $\mathcal{I}$  contains only the pairwise disjoint sets and the total weight of these sets is greater than a given positive number  $\lambda$ , i.e.,  $\sum_{S \in \mathcal{I}} \omega(S) > \lambda$ .

We will show that the weighted  $k$ -set packing problem  $\leq_P$  CRN robust clustering problem. Given an instance of the weighted  $k$ -set packing problem, i.e., a collection of sets  $\mathcal{Q} = \{S_1, S_2, \dots, S_m\}$ , where the set  $S_i, i \in \{1, 2, \dots, m\}$  consists of positive integers. There is an integer weight  $\omega(S_i)$  for  $S_i$ , in the end an integer  $\lambda$  completes the description of this instance. We will construct an instance of a CRN robust clustering problem within polynomial time. W.l.o.g. we let set  $\cup_{i \in \{1, 2, \dots, m\}} S_i = \{1, 2, \dots, N\} = \mathcal{P}$ .

We will construct the CRN and the clusters as follows: For every set  $S \in \mathcal{Q}$ , there will be a corresponding cluster composed with CR nodes constructed. For the set whose size is larger than 1, the IDs of the constructed CR nodes

are identical with the elements in it, and we locate the CR nodes so that any two of them can communicate directly when common channels are available on them. Besides, a set of channels with cardinality of  $|\omega(S)|$  is allocated to all the CR nodes in this cluster, and the channels are on the spectrum band which is exclusive for this cluster. For the set  $S$  which contains only one element, i.e.,  $S = \{t\}$  where  $t \in \mathcal{P}$ , a cluster composed with two CR nodes will be created. In this case, one CR node's ID is  $t$ , the other CR node is the dummy node of the former and its ID is  $t + N$ . A number of  $|\omega(S)|$  channels from the exclusive spectrum band for this cluster are allocated to these two CR nodes. Now we have constructed the clusters which correspond to all the sets in  $\mathcal{Q}$ . Note that every CR node is allowed to form a singleton cluster by itself, although its common channels don't contribute to the sum of  $f(C)$ .

Actually, all the constructed CR nodes can be assumed to locate in a very small area so that each CR node is within the transmission scope of every other CR node. Note that in each constructed cluster, the CR nodes occupy the common channels which are exclusive to this cluster, this design of transformation eliminates the formation of the cluster which doesn't have a corresponding set in  $\mathcal{Q}$ . The existence of the singleton clusters ensures that it is always possible to find out a group of clusters, which together constitute the whole CRN.

Now suppose there is a set of pairwise disjoint clusters which constitute the CRN  $\mathcal{N}$ , and the sum of  $f(C)$  is greater than  $\lambda$ . After removing the singleton clusters, we can easily find the natural association between the remaining clusters and the sets in  $\mathcal{Q}$ . The clusters in the CRN correspond to the sets in  $\mathcal{Q}$  according to the mapping between the node IDs in the clusters and the elements in the sets. In particular, the clusters which contain dummy CR nodes correspond to the sets which contain only one element. Then the sum of the weights of the corresponding sets equals to the sum of  $f(C)$  and thus greater than  $\lambda$ .

We have now shown that our algorithm solves the weighted  $k$ -set packing problem using a black box for the robust clustering problem. Since our construction takes polynomial time, we can conclude that the robust clustering problem is NP-hard.  $\square$