

Versatile Robust Clustering of Ad Hoc Cognitive Radio Network

Di Li, *Member, IEEE*, Erwin Fang, and James Gross, *Member, IEEE*

Abstract—Cluster structure in cognitive radio networks facilitates cooperative spectrum sensing, routing and other functionalities. The availability of unlicensed channels which are available for every member in a cluster decides the robustness of that cluster against the licensed users' influence. Thus there should be more unlicensed channels in the clusters. In the process of forming clusters, every secondary user decides with whom to form a cluster, or which cluster to join. Congestion game model is adopted to analyse this process, which not only contributes the algorithm design directly, but also provides guarantee of convergence into Nash Equilibrium and convergence speed. Our proposed distributed clustering scheme outperforms the comparison scheme in terms of cluster robustness against the primary users, convergence speed and volume of control messages. Furthermore, the proposed clustering solution is versatile to fulfil other requirements such like fast convergence and cluster size control. Besides, we prove the clustering problem is NP-hard, and also propose the centralized solution. The extensive simulation supports our claims.

Index Terms—Cognitive Radio, Cluster, Robust, game theory, congestion game, distributed, centralised, size control.

1 INTRODUCTION

Cognitive radio (CR) is a promising technology to solve the spectrum scarcity problem [1]. Unlicensed users access the spectrum allocated to them whenever there is information to be transmitted. In contrast, unlicensed users can only access the licensed spectrum after validating the channel is unoccupied by licensed users. This refers to the process of sensing a particular channel and verifying (with a previously specified probability of error) it is not used by a primary user currently. In this hierarchical spectrum access model [2], the licensed users are also called primary users (PU), and the CR users are known as secondary users and constitute the cognitive radio networks (CRN).

As to the operation of CRN, efficient spectrum sensing is identified to be critical to the success of cognitive radio networks [3]. Cooperative spectrum sensing is able to effectively cope with noise uncertainty and channel fading, thus remarkably improves the sensing accuracy [4]. Collaborative sensing relies on the consensus of CR users within certain area, and decreases considerably the false sensing reports caused by fading and shadowing of reporting channel. In this regard, clustering is regarded as an effective method in cooperative spectrum sensing [5], [6], as a cluster forms adjacent secondary users as a collectivity to perform spectrum sensing together. Clustering is also efficient to enable all CR

devices within the same cluster to stop payload transmission on the operating channel and initiate the sensing process, so that all the CR users¹ within the one cluster are able to vacate the channel swiftly when primary users are detected by at least one CR node residing in the cluster [7]. With cluster structure, as CR users can be notified by cluster head (CH) or other cluster members about the possible collision, the possibility for them to interfere neighbouring clusters is reduced [8]. Clustering algorithm has also proposed to support routing in cognitive ad-hoc networks [9].

The communication within a cluster is conducted in the spectrum which is available for every member in that cluster. Usually there are multiple unlicensed channels available for all the members in a cluster, which are referred as *common control channels* (CCC). When one or several members can not use one certain CCC because primary users are detected to appear on that channel, this channel will be excluded from the set of CCCs, in particular, if this channel is the working channel, then all the cluster members switch to another channel in the set of CCCs. In the context of CRN, as the activity of primary users is controlled by licensed operators which are generally not known to CR users, the availability of CCCs for the formed clusters is totally decided by primary users' activity. In other words, the availability of CCCs for clusters is passive and can not guaranteed. In CRN, one cluster survives the influence of primary users when at least one CCC is available for that cluster. As the channel occupation by primary users is assumed to be uncontrollable to the CR users, a cluster formed with more CCCs will survive with higher probability. Thus the number of CCCs in one cluster indicates robustness of it when facing ungovernable influence from primary users. As a result, how to form the clusters plays an important role on the robustness of clusters in CRN.

To solely pursue cluster robustness against the primary users' activity, i.e., to achieve more common channels within clusters, the ultimately best clustering strategy is ironically that each node constitutes one single node clusters. Apparently this contradicts our motivation of proposing cluster in cognitive radio network. This contradiction indicates that, the robustness discussed in terms of number of common channels carries little meaning when the sizes of formed clusters are not given consideration. Besides, cluster size plays import roles in certain aspects. For instance, cluster size is one decisive factor in power preservation [10], [11], and it also influences the accuracy of cooperative spectrum

D. Li was with RWTH Aachen university, Germany.

Erwin. Fang is with ETH, Switzerland.

J. Gross is with KTH Royal Institute of Technology, Sweden.

Manuscript received xxxx xx, 20xx; revised xxxx xx, 20xx.

1. The term *user* and *node* are used interchangeably in this paper, in particular, *user* is used when its networking or cognitive ability are discussed or stressed, and *node* is used when the network topology is discussed.

sensing [12]. Hence, cluster size should be given consideration when discussing cluster robustness against primary users.

In this paper, a decentralized clustering approach ROSS (RObust Spectrum Sharing) is proposed to cover the issues of robustness and size control of clusters in CRN. ROSS is able to form clusters with desired sizes, and the generated clusters are more robust than other clustering scheme which has claims on cluster robustness, i.e., more secondary users residing in clusters against increasing influence from primary users. Compared with previous work, ROSS involves much less control messages, and the generated clusters are significantly more robust. We also propose the light weighted versions of ROSS, which involve less overheads and thus are more suitable for mobile networks. Throughout this paper, we refer the clustering schemes on the basis of ROSS as *variants of ROSS*, i.e., the fast versions, or that with size control feature.

The rest of paper is organized as follows. After reviewing related work in section 2, we present our system model in Section 3. Then we introduce our clustering scheme ROSS and its variants in section 4. The clustering problem is given through analysis and a centralized scheme is proposed in section 5. Extensive performance evaluation is in section 6. Finally, we conclude our work and point out direction future research in section 7.

2 RELATED WORK

Prior to the emergence of open spectrum access, as an important method to manage network, clustering has been proposed in for ad hoc networks [13], [14], [15], wireless mesh networks and sensor networks [9]. In ad hoc and mesh networks, the major focus of clustering is to preserve connectivity (under static channel conditions) or to improve routing. In sensor networks, the emphasis of clustering has been on longevity and coverage. Overhead generated by clustering in ad hoc network is analysed in [16], [17].

As to cognitive radio networks, clustering schemes are also proposed, which target different aspects. Work [12] improves spectrum sensing ability by grouping the CR users with potentially best detection performance into the same cluster. Clustering scheme [10] obtains the best cluster size which minimizes power consumption caused by communication within and among clusters. [10] proposes clustering strategy in cognitive radio network, which looks into the relationship between cluster size and power consumption and accordingly controlling the cluster size to decrease power consumption. Cogmesh is proposed in [18] to construct clusters by the neighbour nodes which share local common channels, and by interacting with neighbour clusters, a mesh network in the context of open spectrum sharing is formed. Robustness issue is not considered by this clustering approach. [19] targets on the QoS poisoning and energy efficiency. This approach first decides on the relay nodes which minimize transmission power consumption, then the chosen nodes become cluster heads and clusters are formed in a dynamic coalition process. This work emphasis on power efficiency and doesn't take into account the channel availability and the issue of robustness of the formed clusters. In [6], [20], the channel available to the largest set of one-hop neighbours is selected as common channel which yields a partition of the CRN into clusters. This approach minimizes the set of distinct frequency bands (and hence, the set of clusters) used as common channels within the CRN. However, bigger cluster sizes generally lead to less options within one cluster

to switch to if the common channel is reclaimed by a primary node. Hence, this scheme does not provide robustness to formed clusters. [21] deploys cluster structure in order to implement common channel control, medium access with multiple channel and channel allocation. The node with the maximum number of common channels within its k -hop neighborhood is chosen as cluster head, but how to avoid one node appearing in multiple clusters is not given consideration.

Clustering robustness is considered in [22], [23]. The authors propose a distributed scheme where the metric is the product of cluster size and the number of common control channels. This scheme involves both cluster size and number of CCCs, but it is inherently flawed. With the metric, cluster could be formed only due to one factor of the two, e.g. a spectrum rich node will exclude its neighbour to form a cluster by itself. Besides, this scheme leads to a high variance on the size of clusters, which is not desired in certain applications as discussed in [10], [21].

3 SYSTEM MODEL

We consider a set of cognitive radio users \mathcal{N} and a set of primary users distributed on a two-dimensional Euclidean plane. These users share a number of non-overlapping licensed channels according to the spectrum overlay model. The set of these licensed channels is denoted as \mathcal{K} . As secondary users, the CR users are allowed to transmit on a channel $k \in \mathcal{K}$ only if no primary user is detected being accessing channel k . Further, we consider a *cognitive radio ad-hoc network* which consists of all secondary users and does not contain any primary user.

Secondary users conduct spectrum sensing independently and sequentially on all licensed channels. We assume that every node can detect the presence of primary user on each channel with certain accuracy.² We denote $K_i \subseteq \mathcal{K}$ as the set of available channels for i . We adopt the unit disk model [24] for the transmission of both primary and CR users. Both primary users and CR users have fixed transmission ranges respectively, and all the channels are regarded to be identical in terms of signal propagation. If a CR node locates within the transmission range of primary user p , that CR node is not allowed to use the channel $k(p)$.

We assume that in addition to the licensed channels, there is one dedicated control channel. This control channel could be in ISM band or other reserved spectrum which is exclusively used for transmitting control messages. Actually, the control messages involved in the clustering process can be transmitted on available licensed channels through a rendezvous process by channel hopping [25], [26], i.e., two neighbouring nodes establish communication on the same channel. Over the control channel, a secondary user i can exchange its spectrum sensing result K_i to any $i' \in \text{Nb}(i)$. It is available for any secondary node i to exchange control messages with any other node in its neighborhood $\text{Nb}(i)$ during the cluster formation phase. $\text{Nb}(i)$ is simply defined as the set of nodes located within the transmission range of i .

If a secondary user i is not in the transmission range of a primary user p , i can certainly not detect the presence of p . As the transmission range of primary users is limited and secondary users are located at different locations, different secondary users may have different views of the spectrum availability, i.e., for any $i, i' \in \mathcal{N}$, $K_i = K_{i'}$ does not necessarily hold. As the assumed

2. The spectrum availability can be validated with a certain probability of detection. Spectrum sensing/validation is out of the scope of this paper.

0/1 state of connectivity is solely based on the Euclidean distance between secondary users,

A cognitive radio network can be represented as an undirected graph $G = (N, E)$, where $E \subseteq N \times N$ such that $\{i, i'\} \in E$ if, and only if, there exists a channel $k \in \mathcal{K}$ with $k \in K_i \cap K_{i'}$. Note that we consider the channel availability only for *one* snapshot of time. For the rest of this paper the word channel is referred to licensed channel, if the control channel is not explicitly mentioned.

3.1 Clustering

In this section, we describe what a cluster in the context of CRNs means. A cluster $C \subseteq N$ is a set of secondary nodes consisting of a cluster head h_C and a number of cluster members. The cluster head is able to communicate with any cluster member directly. In other terms, for any cluster member $i \in C$, $i \in \text{Nb}(h_C)$ holds.

Cluster is denoted as $C(i)$ when its cluster head is i . Cluster size of $C(i)$ is written as $|C(i)|$. $K(C) = \cap_{i \in C} K_i$, $K(C)$ denotes the set of common control channels in cluster C . Clustering is performed periodically, as secondary users are mobile and the channel availability on secondary users are changing as primary users change their operation state.

TABLE 1
Notations in robust clustering problem

| Symbol | Description |
|----------------|---|
| N | collection of secondary users, $N = 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 CCCs of cluster $C(i)$ |
| h_C | the cluster head of a cluster C |
| δ | desired cluster size |
| 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 | social connectivity degree of CR node i |
| C_i | the i th legitimate cluster (only appear in Sec. 5.2) |

4 CLUSTERING ALGORITHM: ROSS

In this section we introduce the distributed clustering scheme ROSS which leads to robust clusters against PU' influence. ROSS consists of two cascaded phases: *cluster formation* and *membership clarification*. With ROSS, CR nodes form clusters on the basis of the proximity of the available spectrum in their neighbourhood. Afterwards, CR nodes belong to one certain cluster.

4.1 Phase I - Cluster Formation

After conducting spectrum sensing and communication with neighbours, every CR node is aware of the available channels available for themselves as well as for all its neighbors. For each CR user i , two metrics are proposed to characterize the spectrum proximity between u and its neighborhood:

- *Individual connectivity degree* d_i : $d_i = \sum_{j \in \text{Nb}(i)} |K_i \cap K_j|$. It denotes the sum of the numbers of common controls channels

between node i and every neighbour. It is an indicator of node i 's adhesive property to the CRN.

- *Social connectivity degree* g_i : $g_i = |\bigcap_{j \in \text{Nb}(i) \cup i} K_j|$. It is the number of common channels available to i and all its neighbors. g_i represents the ability of i to form a robust cluster with its neighbours.

Individual connectivity degree d_i and social connectivity degree g_i together form the *connectivity vector*. Figure 1 illustrates an example CRN where each node's connectivity vector is calculated and shown.

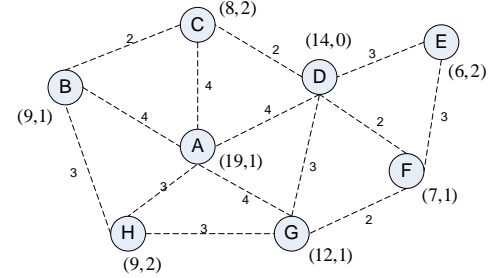


Fig. 1. Connectivity graph and the connectivity vector (d_i, g_i) for each node. The available channels sensed by each node are: $K_A = \{1, 2, 3, 4, 5, 6, 10\}$, $K_B = \{1, 2, 3, 5, 7\}$, $K_C = \{1, 3, 4, 10\}$, $K_D = \{1, 2, 3, 5\}$, $K_E = \{2, 3, 5, 7\}$, $K_F = \{2, 4, 5, 6, 7\}$, $K_G = \{1, 2, 3, 4, 8\}$, $K_H = \{1, 2, 5, 8\}$. Dashed edge indicates the end nodes are within each other's transmission range. The number of common channels between two nodes is shown by the edge.

After introducing the connectivity vector, we proceed to introduce the first phase of algorithm ROSS. To put it briefly, in this phase cluster heads are determined in the beginning. Then clusters are formed on the basis of the cluster heads' neighborhoods.

4.1.1 Determining Cluster Heads and Form the Initial Clusters

In this phase, 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 any neighbors except for those which have already become cluster heads (the appearance of cluster heads will be explained in Section 4.2), then node i becomes clusters 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 $d_j < d_k, \forall k \in \text{Nb}(j) \setminus \{\text{CHs} \cup i\}$, then the node out of $\{i, j\}$ with higher social connectivity degree becomes cluster head. The other nodes become a member of that cluster. 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's head. The node which becomes cluster head broadcasts a message of its eligibility of being cluster head to notify its neighbours, and claims its neighbourhood as its cluster. The pseudo code for the cluster head decision and the initial cluster formation is shown in Algorithm 1 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 $M > |\mathcal{K}| \cdot N$. Then i broadcasts its new individual connectivity degree to all its neighbors. We manipulate the individual connectivity degree of the CR nodes which are included in certain clusters. Hence, nodes located outside of the a formed cluster can possibly become cluster heads or can also be included into other clusters. When a CR node i is associated to multiple clusters, i.e., i has received multiple notifications of cluster head eligibility from different CR nodes, d_i is still set to M . We have the following theorem to show

that as long as a secondary user's individual connectivity degree is greater than zero, that secondary user will eventually be integrated into a certain cluster, or it eventually becomes a cluster head.

THEOREM 4.1: *Given a CRN, every secondary user is included into at least one cluster within N steps. (COMMENT: N has to be defined within the theorem!!)*

Here, by *step* we mean one application of Algorithm 4.1 by a secondary user. The Proof is in Appendix 19. According to Theorem 4.1, Phase I completes within a reasonable amount of time.

Let us apply Algorithm 1 to the example shown in Figure 1. 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 Guarantee the Availability of Common Control Channel

After executing phase I of ROSS, there are some secondary users which become cluster heads. The cluster head and its neighbourhood (except for the CHs) form a cluster. It is possible that certain formed clusters don't own CCCs, then we solve this problem with the following method.

As decreasing cluster size increases the number of CCCs within the cluster, to have at least one CCC, certain nodes are eliminated. 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 number of 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 its neighbours and resulting in a *singleton cluster* which is composed by itself. The pseudo code for cluster head to obtain at least one common channel is shown in Algorithm 2. As to the nodes eliminated in this procedure, they restore their original individual connectivity degrees, and become either cluster heads or get included into other clusters afterwards according to Theorem 4.1.

4.1.3 Cluster Size Control in Dense CRN

In the introduction section, we have stated that cluster size should be given consideration to justify the concept of robustness of clusters, i.e., without specifying requirement on cluster sizes, small clusters will be generated to obtain more CCCs. Except for cooperative sensing, clusters need to conduct some other functionalities. When cluster size is large, there will be substantial burden on cluster heads to manage the cluster members, which is a challenge for resource limited cluster heads, thus the cluster size should fall in a desired range [27], [28].

In the following we illustrate the pressing necessity to control the cluster size when CRN becomes dense via both theoretical analysis and simulation. Assuming the secondary and primary users are evenly distributed and primary users occupy the licensed channels randomly, then both CR nodes density and channel availability in the CRN can be seen to be spatially homogeneous. The formed clusters are the neighbourhoods of cluster heads, and the neighbourhood is decided by the transmission range and network density. We consider a cluster $C(i)$ where i is CH in a dense CRN. When we don't consider the CHs which could appear within i 's

neighbourhood in the procedure of guaranteeing CCCs, according to Algorithm 1, the nearest cluster heads could locate just outside node i 's transmission range. An instance of this situation is shown in Figure 2. In the figure, black dots represent cluster heads, the circles denotes the transmission ranges of cluster heads. Cluster members are not shown in the figure. Let l be the length of side of

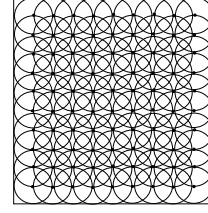


Fig. 2. Clusters formation in extremely dense CRN. Black dots are cluster heads, other cluster members are not drawn.

simulation plan square, and r be CR's transmission radius. Based on the aforementioned analysis and geometry illustration as shown in Figure 2, we give an estimate on the maximum number of generated clusters, which is the product of the number of cluster heads in one row and that number in one line, $l/r * l/r = l^2/r^2$. Given $r=10$ and $l=50$, the number of formed clusters is shown in Figure 3. With the increase of CR users, network density increases linearly (the Y axis label is the average number of neighbours), and the number of formed clusters also increases and approaches to the the upper bound of 25 which complies with the estimation.

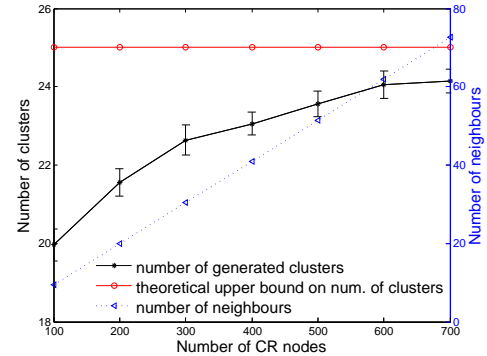


Fig. 3. The correlation between the number of formed clusters and network density. Note that the number of neighbours denotes the network density. Simulation is run for 50 times and the confidence interval is 95%.

Both the analysis and simulation show that when applying ROSS, after the number of clusters saturates with the increase of network density, the cluster size increases almost linearly with the network density, thus certain measures are needed to curb this problem. This task falls to the cluster heads.

To control the cluster size, cluster heads prune their cluster members to achieve the desired cluster size. The desired size δ is decided based on the capability of the CR users and the tasks to be conveyed. Given the desired size δ , a cluster head excludes members sequentially according to the following principle, the absence of one cluster member leads to the maximum increase of common channels within the cluster. This process ends when the size of resultant cluster is δ and at least one CCC is available. This procedure is similar with that to guarantee CCCs in cluster, thus the algorithm can reuse Algorithm 2.

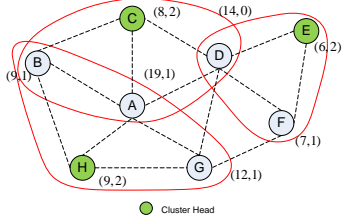


Fig. 4. Clusters formation after the first phase of ROSS. There are some nodes being debatable nodes, i.e., belonging to more than one cluster.

4.2 Phase II - Membership Clarification

After applying phase I of ROSS to the example in Figure 1, the resulted clusters are shown in Figure 4. We notice nodes A, B, D are included in more than one cluster. We refer these nodes as *debatable nodes* as their cluster affiliations are not clear, and the clusters which include the debatable node i are called *claiming clusters* of node i , and are denoted as S_i . Actually, debatable nodes extensively exist in CRN with larger scale. Debatable nodes should be exclusively associated with only one cluster and be removed from the other claiming clusters, this procedure is called *cluster membership clarification*. We will introduce the solution for cluster membership clarification in the following.

4.2.1 Distributed Greedy Algorithm (DGA)

After Phase I, debatable nodes, e.g., i needs to decide one cluster $C \in S_i$ to stay, and thereafter leaves the rest others in S_i . The principle for debatable node i to choose one claiming cluster is that its decision can result in the greatest increase of common channels in all its claiming clusters. Since node i is a neighbour of all the cluster heads in S_i , node i is aware of the channel availability on these claiming cluster heads, and the common control channels in these claiming clusters. With these information, node i is able to calculate how many more CCCs will be produced in one claiming cluster if i leaves that cluster. If there exists one cluster $C \in S_i$, when i leaves this cluster brings the least increased CCCs than leaving any other claiming clusters, then i chooses to stay in cluster C . When there comes a tie, among the relevant claiming clusters, i chooses to stay in the cluster whose cluster head shares the most CCCs with i . In case there are multiple claiming clusters demonstrating the same on the aforementioned criteria, 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 the previous rules could not decide on the unique claiming cluster to stay. The pseudo code of this algorithm is described as Algorithm 3. After deciding its membership, debatable node i notifies all its claiming clusters, and retrieves the updated information of the claiming clusters, e.g., $K(C)$, K_{hc} , $|C|$, where $C \in S_i$.

This procedure raises a concern that the debatable nodes may never stop changing their affiliations, as debatable nodes' choices seem to be dependent on each other, thus the infinite chain effect never ceases. For example, assuming one debatable node i locates in cluster $C \in S_i$, and C has more than one debatable node except for i . Assuming node i makes decision to stay in the claiming cluster C , afterwards one another debatable nodes j decides its affiliation, and there is $j \in C \in S_i$. When j leaves cluster C , which decrease cluster C 's cluster size, and could possibly trigger node i to alter its previous decision to leave C , as C 's size is smaller now and leaving it may result in more increase of CCCs in S_i . At this point of time, debatable node j may rejoin cluster

C due to the changes in S_j , then both node i and j face the *same*³ situation again. Thence, we need to answer this concern before implementing ROSS-DGA. In the following we show that the process of membership clarification can be formulated into a game, and a 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, modelling and analysing the interactions among individuals. A game consists of three elements: a set of players, a selfish utility for each player, and a set of 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 [29]. In the past few years, game theory has been extensively applied to problems in communication and networking [30], [31]. Congestion game is an attractive game model which describes the problem where participants compete for limited resources in a non-cooperative manner, it has good property that Nash equilibrium can be achieved after finite steps of best response dynamic, i.e., each player choose strategy to maximizes/minimizes its utility/cost with respect to the other players' strategies. 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 [32], or users downloading files from cloud, etc.

To formulate the debatable nodes' membership clarification into the desired congestion game, we observe this process from a different (or opposite) perspective. From the new perspective, the debatable nodes are regarded to be isolated and don't belong to any cluster, in other words, their claiming clusters become clusters which are beside them. Now for the debatable nodes, the previous problem of deciding which clusters to leave becomes a new problem that which cluster to join. In the new problem, debatable node i chooses one cluster C out of S_i to join if the decrease of CCCs in cluster C is the smallest in S_i , and the decrease of CCCs in cluster C is $\sum_{C \in S_i} \Delta|K(C)| = \sum_{C \in S_i} (|K(C)| - |K(C \cup i)|)$. The interaction between the debatable nodes and the claiming clusters is shown in Figure 5.

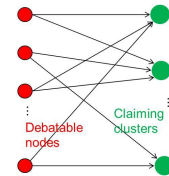


Fig. 5. Debatable nodes and claiming clusters

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 $\Gamma = (\mathcal{P}, \mathcal{R}, \sum_{i \in \mathcal{P}}, f)$, and the elements in Γ are explained below,

3. Actually it is not totally same as before, as there are some new changes within S_j .

- \mathcal{P} , the set of players in the game, which are the debatable nodes in our problem.
- $\mathcal{R} = \cup_{i \in \mathcal{P}} S_i$, denotes the set of resources for players to choose, in our problem, S_i is the set of claiming clusters of node i , and \mathcal{R} is the set of all claiming clusters.
- Strategy space $\sum_i, i \in \mathcal{P}$ is the set of 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 decrease of CCCs in cluster C when debatable node i joins C . As to cluster $C \in S_i$, the decrease of CCCs caused by the enrolment of 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 [33]. 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 CCCs $\sum_{i: C \in S_i} \Delta|K^i(C)|$, 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 of number of steps before convergence is $n^2 * m$ [33], 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 upper bound of the number of steps can be expressed as $O(N^3)$.

In fact, the number of steps which are actually involved in this process is much smaller than N^3 , as both n and m are considerably smaller than N . The percentage of debatable nodes in \mathcal{N} is illustrated in Figure 14, which is between 10% to 60% of the total number of CR nodes in the network. The number of clusters heads, as discussed in Section 4.1, is dependent on the network density and the CR node's transmission range. As shown in Figure 3, the cluster heads take up only 3.4% to 20% of the total number of CR nodes.

4.2.3 Distributed Fast Algorithm (DFA)

We propose a faster version of ROSS, ROSS-DFA, which differs from ROSS-DGA in the second phase. With ROSS-DFA, debatable nodes decide their respective cluster heads 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 decision 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 dynamically and re-clustering is necessary. To run ROSS-DFA, debatable node executes only one loop in Algorithm 3.

Now we apply both ROSS-DGA and ROSS-DFA to the toy network in Figure 4 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 6. In Figure 6(a), node A staying in $C(C)$ and leaving $C(H)$ brings 2 more CCCs to S_A , which is more than that brought by another strategy showed in 6(b). After the decisions made similarly by the other debatable nodes B and D , the final clusters are formed as shown in Figure 7.

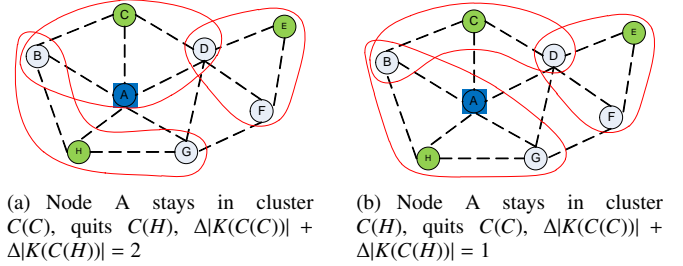


Fig. 6. Membership clarification: possible cluster formations caused by node A's different choices

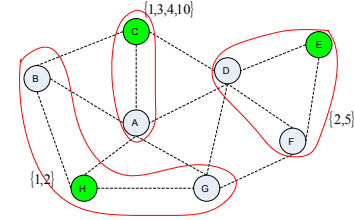


Fig. 7. Final formation of clusters, CCCs for each cluster is shown. $K(C(C))$, $K(C(E))$, $K(C(H))$ are shown beside corresponding clusters.

5 CENTRALIZED CLUSTERING SCHEME

The centralized clustering scheme aims to form clusters with desired sizes, meanwhile the total number of CCCs of all clusters is maximized. In the following, we refer this problem as *centralized clustering*, and give the formal problem definition.

DEFINITION 1: Centralized clustering in CRN.

Given a cognitive radio network \mathcal{N} where nodes are indexed from 1 to N sequentially. Based on certain correlation, certain secondary users constitute one cluster C . $1 \leq |C| \leq k$ where $|C|$ is the size of cluster C and k is a positive integer. We name the collection of such clusters as $\mathcal{S} = \{C_1, C_2, \dots, C_{|\mathcal{S}|}\}$, where \mathcal{S} satisfies the following properties: $\bigcup_{1 \leq i \leq |\mathcal{S}|} C_i = \mathcal{N}$ and $K(C(i)) \neq \emptyset$ for any i which satisfies $1 \leq i \leq |\mathcal{S}|$.

We give a new definition of the number of CCCs, where the number of common control channels is $|K(C)|$ if $|C| > 1$, and is zero when $|C| = 1$. We use $f(C)$ to denote the number of CCCs of a cluster C in the new definition.

The centralized clustering problem is to find a subcollection $\mathcal{S}' \subseteq \mathcal{S}$, so that $\bigcup_{C_j \in \mathcal{S}'} C_j = \mathcal{N}$, and $C_j \cap C_{j'} = \emptyset$ for $C_j, C_{j'} \in \mathcal{S}'$ and $j' \neq j$, so that $\sum_{C \in \mathcal{S}'} f(C)$ is maximized. The decision version of centralized clustering in CRN is to ask whether there exists a non-empty $\mathcal{S}' \subseteq \mathcal{S}$, so that $\sum_{C \in \mathcal{S}'} f \geq \lambda$ where λ is a real number.

5.1 Complexity of Clustering Problem

In this section we investigate the complexity of centralized clustering problem. Theorem 5.1 tells centralized clustering problem in CRN is one NP-hard problem.

THEOREM 5.1: *CRN clustering problem is NP-hard, when the maximum size of clusters $k \geq 3$.*

The proof is in Appendix 19.

5.2 Centralized Optimization

As there is no efficient algorithm to solve clustering problem in CRN, we propose a centralized optimization where the objective function and the constraints are heuristic, then we adopt binary linear programming to solve the problem.

Given a CRN \mathcal{N} and desired cluster size δ , we obtain a collection of clusters \mathcal{G} which contains all the *legitimate* clusters, and the sizes of these clusters are $1, 2, \dots, \delta$. Legitimate clusters are the clusters which satisfy the conditions in Section 3.1. Note that the legitimate clusters include the singleton ones, so that we can guarantee the partition of any network is always feasible.

With $N = |\mathcal{N}|$, $G = |\mathcal{G}|$, we construct a constant $G \times N$ matrix $Q_{G \times N}$. The element of matrix Q is q_{ij} , where the subscript i is the index of legitimate cluster, and j is the node ID of one CR node. There are $i \in \{1, 2, \dots, G-1, G\}$, and $j \in \{1, 2, \dots, N-1, N\}$. Element $q_{ij} = |K(C_i)|$ if node $j \in C_i$, and $q_{ij} = 0$ if $j \notin C_i$. In other words, each non-zero element q_{ij} denotes the number of CCCs of the cluster i where node j resides.

$$\begin{matrix}
 & 1 & 2 & 3 & \cdots & j & \cdots & N-1 & N \\
 \begin{matrix} 1 \\ 2 \\ \vdots \\ i \\ \vdots \\ \vdots \\ G \end{matrix} & \begin{pmatrix} |K(C_1)| & |K(C_1)| & 0 & \cdots & \cdots & \cdots & 0 & 0 \\ |K(C_2)| & 0 & |K(C_2)| & \cdots & \cdots & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & |K(C_i)| & 0 & \cdots & \cdots & \cdots & |K(C_i)| & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & 0 & 0 & \cdots & \cdots & \cdots & |K(C'_i)| & 0 \\ |K(C_G)| & \cdots & \vdots & \cdots & \cdots & \vdots & \vdots & \vdots \end{pmatrix}
 \end{matrix}$$

Fig. 8. An example of Matrix Q , its rows correspond to all legitimate clusters, and columns correspond to the CR nodes in the CRN.

We build a $G \times N$ binary variable matrix X , which illustrates the clustering solution. The element of matrix X is binary variable x_{ij} , $i = 1, \dots, G$, $j = 1, \dots, N$. Now, we can formulate the optimization problem as follows,

$$\begin{aligned}
 \min_{x_{ij}} \quad & \sum_{j=1}^N \sum_{i=1}^G (-x_{ij} q_{ij} + (1 - w_i) * p) \\
 \text{subject to} \quad & \sum_{i=1}^G x_{ij} = 1, \text{ for } \forall j = 1, \dots, N \\
 & \sum_{j=1}^N x_{ij} = |C_i| * (1 - w_i), \text{ for } \forall i = 1, \dots, G \\
 & x_{ij} \text{ and } w_i \text{ are binary variables.} \\
 & i \in \{1, 2, \dots, G\}, \quad j \in \{1, 2, \dots, N\}
 \end{aligned}$$

The objective is a sum of two parts, the first part is the sum of products of cluster size and the corresponding number of CCCs, which is the sole metric adopted by the scheme SOC [22]. The second part is the *punishment* for choosing the clusters whose sizes are not δ . In fact, the second part is particularly designed to eliminate the drawbacks of SOC, i.e., SOC produces a large number of singleton clusters and a few large clusters. In practical computation, we minimize the opposite of the first part, and the punishment is a positive value. The first constraint restricts each node j to reside in exactly one cluster. In the second constraint, w_i

is an auxiliary binary variable which denotes whether cluster C_i is chosen by the solution, in particular,

$$w_i = \begin{cases} 0 & \text{if } i\text{th legitimate cluster } C_i \text{ is chosen} \\ 1 & \text{if } i\text{th legitimate cluster } C_i \text{ is not chosen} \end{cases}$$

The second constraint regulates when the i th legitimate cluster C_i is chosen, the number of elements which are 1 in the i th row of the matrix X is $|C_i|$. Now we explain how does the mechanism of the punishment in the objective work. The parameter p is defined as follows,

$$p = \begin{cases} 0 & \text{if } |C_i| = \delta \\ \alpha_1 & \text{if } |C_i| = \delta - 1 \\ \alpha_2 & \text{if } |C_i| = \delta - 2 \\ \vdots & \end{cases}$$

where $\alpha_i > 0$ and increases when $|C_i|$ diverges from δ . Because of w_i , any chosen cluster ($w = 0$) brings certain *punishment*. When the chosen cluster's size is desired size δ , the punishment is zero. In contrary, when the chosen cluster's size diverges from δ , the objective function suffers *loss*. In particular, when $w_i = 0$ and $|C_i| = 1$, the punishment is the most severe. This design doesn't follow the definition of $f(C)$ in Definition 1 strictly, where $f(C_i) = 0$ when $|C_i| = 1$, but our design echoes the definition by exerting the most severe punishment on the singleton clusters in the clustering solution. Choice of α_i affects the resultant clusters.

The optimization formulation is an integer linear optimization problem, which is solved by the function *bintprog* provided in MATLAB. Note that the proposed centralized solution is heuristic. We reiterate the reasons for pursuing the heuristic scheme, first, the problem of centralized clustering is NP hard, and there is no efficient solution to solve it. The second reason is, the collection of legitimate clusters is dependant on the network topology and spectrum availability in the network, thus to each specific CRN, the space of solution is different.

5.2.1 Example of the Centralized Optimization

We look into how does the centralized scheme perform in the toy example of the CRN in Figure 1. We let the desired cluster size δ be 3. A collection of clusters \mathcal{G} is obtained, which contains all the clusters satisfying the conditions of cluster in Section 3.1 and the sizes of clusters are 1, 2 or 3. $\mathcal{G} = \{\{A\}, \{B\}, \dots, \{B, C\}, \{B, A\}, \{B, H\}, \dots, \{B, A, C\}, \{B, H, C\}, \{A, D, C\}, \dots\}$, and $G = |\mathcal{G}| = 38$.

When α_1 and α_2 are set as 0.2 and 0.8, the formed clusters are shown in Fig. 9(b). The resulted clustering solutions from ROSS-DGA/DFA and SOC are shown in Fig. 7 and Fig. 9(a) respectively.

As to the average number of CCCs, the results of ROSS (including both ROSS-DGA and ROSS-DFA), centralized and SOC are 2.66, 2.66, and 3 respectively. Note there is one singleton cluster $C(H)$ generated by SOC, which is not preferred. When we take no account of the singleton clusters, then the average number of common channels of SOC drops to 2.5.

6 PERFORMANCE EVALUATION

In this section, we evaluate the performances of all the variants of ROSS, i.e., ROSS-DGA and ROSS-DFA, and that with cluster size control features. The latter is referred as ROSS- δ -DGA/ROSS- δ -DFA, where x is the desired cluster size. We choose SOC as comparison scheme. To the best of our knowledge, SOC [22] is the only work emphasizing on the robustness of clustering structure from all previous work on clustering in CRN. The authors of [22]

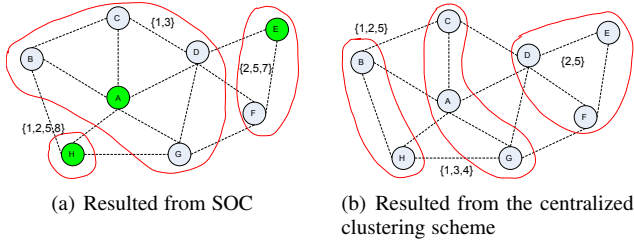


Fig. 9. Final clusters formed in the example network when being applied with SOC and the centralized clustering scheme.

compared SOC with other schemes in terms of the average number of CCCs 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 [15], or for CRN but just considering connection among CR nodes [6]. Hence, we only compare our proposed schemes with SOC to show the ROSS's merits as a distributed scheme, and compare with the centralized scheme to examine the gap with the global optima. In particular, we investigate the following metrics.

- *Average number of CCCs per non-singleton cluster.* This metric shows the robustness of the current non-singleton clusters. Non-singleton cluster refers the cluster whose cluster size is larger than 1. SOC [22] compares the average number of CCCs per cluster without distinguishing singleton clusters, which is biased as singleton clusters don't contribute to the cluster structure.
- *Number of singleton clusters.* This is a straight forward metric which reflects the effectiveness of clustering scheme. The less resulted singleton clusters means more secondary users are benefited from cluster structure. When we investigate the performance with moderate and vigorous intensity of primary users' activities, this metric is the antonym of *survival rate*, i.e., the percentage of nodes which are still within certain clusters.
- *Cluster sizes.* Specific clusters size is pursued in many applications due to energy preservation and the system design [10]. We will present 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.

The simulation is written in C++. Certain number of 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%. All primary and CR users are assumed to be static during the process of clustering. Simulation consists of two parts, in the first part, we investigate the performance of centralized scheme, and the gap between the distributed schemes and the centralized scheme. This part is conducted 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.

6.1 Centralized Schemes vs. Decentralized Schemes

10 primary users and 20 CR users are dropped randomly (with uniform distribution) within a square area of size A^2 , where we set the transmission ranges of primary and CR users to $A/3$. With this setting, the average number of neighbours of one CR user is 4.8. CR users are assumed to be able to sense the existence of primary

users and identify available channels. When clustering scheme is executed, around 7 channels are available on each CR node.

The desired cluster size δ is 3, the parameters used in the *punishment* for choosing the clusters with undesired sizes are set as follows, $\alpha_1 = 0.4$, $\alpha_2 = 0.6$. Performance results are averaged over 50 randomly generated topologies, and the confidence interval corresponds to 95% confidence level.

6.1.1 Number of CCCs in Non-singleton Clusters

Figure 10 shows the average number of the CCCs of non-singleton clusters, from which we can see the centralized schemes outperform the distributed schemes. As to the distributed schemes, SOC achieves the largest number of CCCs than all the variants of ROSS. The reason is, SOC is liable to group the neighbouring CRs which share the most abundant spectrum together, no matter how many of them are, thus the number of CCC of the formed clusters is higher. But this method leaves considerable number of CRs to form singleton clusters.

As to the variants of ROSS, greedy mechanism increases CCCs in non-singleton clusters significantly. We also notice that the size control feature doesn't affect the number of CCCs for both ROSS-DGA and ROSS-DFA. Size control mechanism converts the large clusters into small ones, but meanwhile clusters with the desired sizes have to be made when forming smaller clusters is possible.

6.1.2 Number of Singleton Clusters

Initially, clusters are formed with the presence of 10 PUs. Afterwards, extra 20 batches of PUs are added to the network sequentially, and each batch includes 5 PUs. Figure 11 shows the number of unclustered CRs with the increasing number of PUs, from which two conclusions are drawn.

- The centralized scheme with desired size of two generates the most robust clusters, and SOC results in the most vulnerable clusters. The centralized scheme with desired size of 3 doesn't surpass the variants of ROSS. The reason lies in the the uniformly sized clusters, whereas the variants of ROSS generate considerable amount of smaller clusters which are more likely to survive when PUs' activities become intense. ROSS with size control distinguishes itself as the size control avoid the appearance of the clusters with large size.
- Greedy algorithm improves the survival rate. Greedy strategy adopted in the second phase of ROSS improves the robustness of clusters, i.e. ROSS-DGA exceeds ROSS-DFA, and ROSS- δ -DGA surpasses ROSS- δ -DFA. When the debatable CRs greedily update their affiliations, one of the metrics is the maximum increase of CCCs of the demanding clusters. This observation complies with the result shown in Figure 10.

6.1.3 Cluster Size Control

Figure 12 depicts the empirical cumulative distribution of the CRs residing in certain sized clusters which are generated in 50 runs. The centralized schemes form clusters which satisfy the requirement on cluster size strictly. As to ROSS-DGA and ROSS-DFA with size control mechanism, CR nodes averagely distributes 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 is 4.5, 5, and 5.5, and the 90% percentiles for the three schemes are 8, 8, and 9.

Note ROSS-DGA and ROSS-DFA with size control feature generate 10%-20% singleton clusters, which is due to the cluster



Fig. 10. Number of common channels for non-singleton clusters

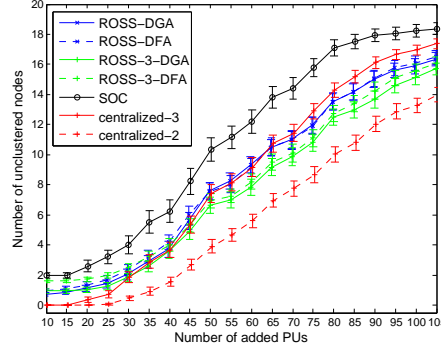


Fig. 11. Number of CRs which are not included in any clusters

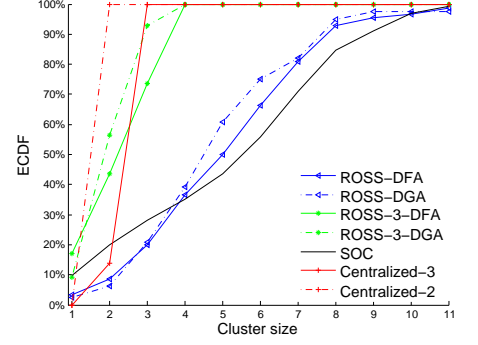


Fig. 12. Cumulative distribution of CRs residing in clusters with different sizes

Fig. 13. Comparison between the distributed and centralized clustering schemes in a small network ($N = 20$)

pruning discussed in section 4.1.3, whereas, without size control, only 3% nodes are in singleton clusters. When applying SOC, 10% of nodes are in singleton clusters.

6.1.4 Control Signalling Overhead

In this section we compare the amount of control messages involved in different clustering schemes, i.e., centralized scheme, ROC, and the variants of ROSS. We omit the control messages involved in the process of neighbourhood discovery, which is the premise for any clustering scheme. According to [34], the message complexity is defined as the number of messages used by all nodes. To have the same criterion to compare the overhead of signalling, we count the number of transmissions of control messages, without distinguishing they are sent with broadcast or unicast. This metric is synonymous with the number of updates discussed in Section 4.

As to ROSS, the control messages are generated in both phases. In the first phase, when a CR node decides itself to be the cluster head, it broadcasts a message containing its ID, cluster members and the set of CCCs in its cluster. In the second phase, a debatable node broadcasts its affiliation to inform its claiming clusters, then the CHs of the claiming clusters broadcast message about the new cluster members if they are changed due to the debatable node's decision. The total number of the decisions involved in cluster formation has been analysed in Theorem 4.1 and Section 4.2.2 respectively.

Comparison scheme SOC involves three rounds of execution. In the first two rounds, every CR node maintains its own cluster and seek to integrate neighbouring clusters, or joins one of them. The final clusters are obtained in the third round. In each round, every CR node is involved in comparisons and cluster mergers.

The centralized scheme is conducted at the centralized control device, but it involves two phases of control message transmission. The first phase is information aggregation, in which every CR node's channel availability and neighborhood is transmitted to the centralized controller. The second phase is broadcasting, where the clustering solution is disseminated to every CR node.

We adopt the algorithm proposed in [35] to broadcast and gather information as the algorithm is simple and self-stabilizing. This scheme needs building a backbone structure to support the communication, we use our generated cluster heads as the backbone and the debatable nodes as the gateway nodes between

the backbone nodes. As the backbone is built once and can support transmission for multiple times, the messages involved in the clustering process are not counted. As to the process of information gathering, we assume that every cluster member sends the spectrum availability and its ID to its cluster head, which further forwards the message to the controller, thus the number of transmission is N . As to the process of dissemination, in an extreme situation where all the gateway and the backbone nodes broadcast, the number of transmission will be $h + m$, where h is the number of cluster heads, m is number of debatable nodes, d is the average number of demanding clusters for each debatable node.

The number of control messages which are involved in both ROSS and the centralized scheme is related with the number of debatable nodes. Figure 14 shows the percentage of debatable nodes when the CRN network becomes denser, from which we can obtain the value of m .

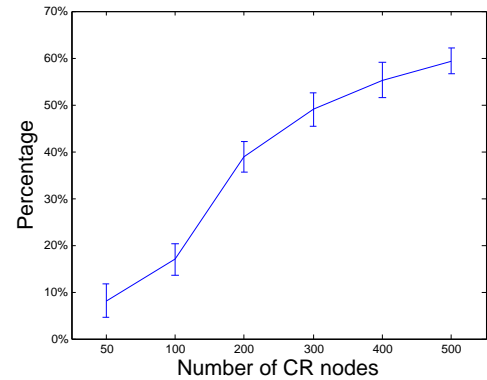


Fig. 14. The percentage of debatable nodes after phase I of ROSS.

The message complexity, quantitative analysis of the number of control messages involved in clustering, and the size of control messages are shown in Table 2. Figure 15 shows the number of transmissions of SOC, the upper bound of the number of transmissions for ROSS, and the analytical number of transmissions of the centralised scheme.

TABLE 2
Singalling overhead

| Scheme | Message Complexity | Quantitative number of messages | Content of message (size of message) |
|----------------------------------|-----------------------|---------------------------------|---|
| ROSS-DGA, ROSS- δ -DGA | $O(N^3)$ (worst case) | $h + 2 * m^2 d$ (upper bound) | Phase I: notification from cluster head (1 byte), new individual connectivity degree (1 byte); Phase II: update of debatable nodes' affiliation (1 byte), claiming cluster i ' new membership ($ C(i) $ bytes) |
| ROSS-DFA, ROSS- δ -DFA | $O(N)$ (worst case) | $h + 2m$ (upper bound) | |
| SOC | $O(N)$ | $3 * N$ | $C(i)$ ($ C(i) $ bytes), $K(C(i))$ ($ P $ bytes), $i \in N$ |
| Centralized | $O(N)$ | $h + m + N$ (upper bound) [35] | $\{C\}$ ($ C_i * N$ bytes) |

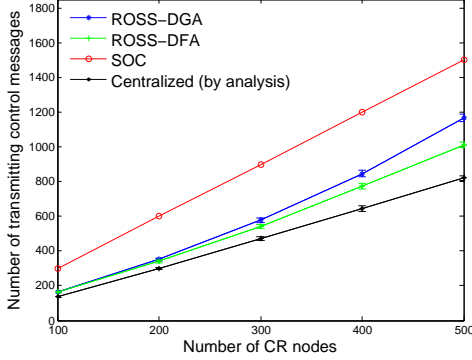


Fig. 15. Number of control messages. Note the curves for ROSS-DGA and ROSS-DFA are the upper bounds of the number of messages, the curve of centralized scheme reflects an extreme situation.

6.2 Comparison between Distributed Schemes

In this section we investigate the performances of distributed clustering schemes in CRN with different network scales and densities. The transmission range of CR is $A/10$, PR's transmission range is $A/5$. The initial number of PU is 30. The desired sizes adopted are listed in the Table 3.

TABLE 3

| | | | |
|----------------------------|-----|-----|-----|
| Number of CRs | 100 | 200 | 300 |
| Average num. of neighbours | 9.5 | 20 | 31 |
| Desired size δ | 6 | 12 | 20 |

6.2.1 Number of CCCs per Non-singleton Clusters

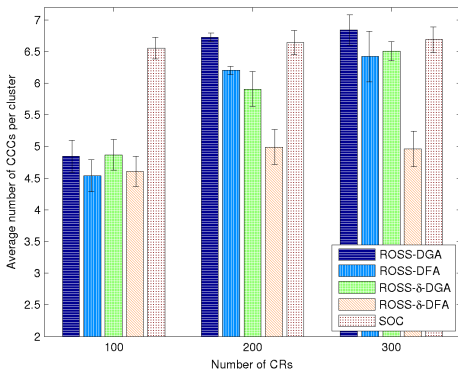


Fig. 16. Number of common channels of non-singleton clusters.

Figure 16 illustrates the average number of CCCs of the non-singleton clusters. When $N = 100$, variants of ROSS have 30% less CCCs than SOC, but this gap is decreased significantly when N is 200 and 300, i.e., when $N = 300$, the number of CCCs achieved by ROSS variants (except for ROSS- δ -DFA) is almost the same as that resulted from SOC. This means SOC outperforms in terms of the average number of CCCs per non-singleton cluster when network is sparse.

this is also observed in the evaluation in Section 6.1.1 where $N = 20$. When the network becomes denser, even this metric favours SOC as discussed in the beginning of Section 6, ROSS-DGA achieves even more CCCs than SOC, and ROSS-DFA and ROSS- δ -DGA increase the number of CCCs visibly.

6.2.2 Number of Singleton Clusters

Figure 17 illustrates the increasing trend of singleton clusters with the increase of PUs, when $N = 100$. SOC generates around 10 more singleton clusters than the variants of ROSS, which accounts for 10% of the total CR nodes. We only show the average values of the variants of ROSS as their confidence intervals overlap. Figure 18 depicts a denser CRN where $N = 300$. SOC noticeably causes more singleton clusters than ROSS variants, except that ROSS-20-DFA results in more singleton clusters when PUs are few. The reason is ROSS-20-DFA conducts cluster membership clarification for only once, which causes large number of singleton clusters. ROSS-20-DGA increases the size of smaller clusters through debatable nodes' repeated updates thus drastically decreases the number of singleton clusters.

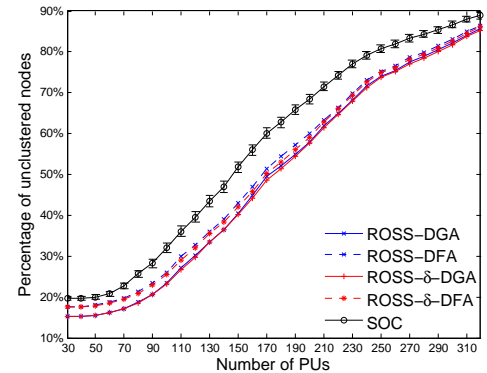


Fig. 17. Percentage of CRs which are not included in any clusters with the increasing number of primary users, $N = 100$

From the Figure 17 and 18, we can conclude that the clusters obtained from the variants of ROSS are more robust than SOC. Besides, the greedy mechanism moderately strengthens the robustness of the clusters.

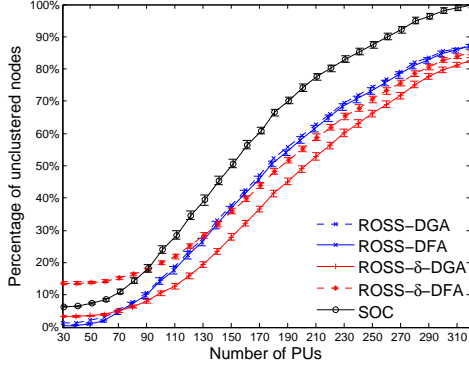


Fig. 18. Percentage of CRs which are not included in any clusters with the increasing number of primary users, $N = 300$

6.2.3 Cluster Size Control

As Fig. 19 shows, when the network scales up, the number of formed clusters by ROSS increases by smaller margin. This result coincides with the analysis in Section 4.1.3, that the number of formed clusters saturates when the network scales. When the network becomes denser, more clusters are generated by SOC compared with ROSS variants. To better understand the distribution of the sizes of formed clusters, we depict the cluster sizes with cumulative distribution.

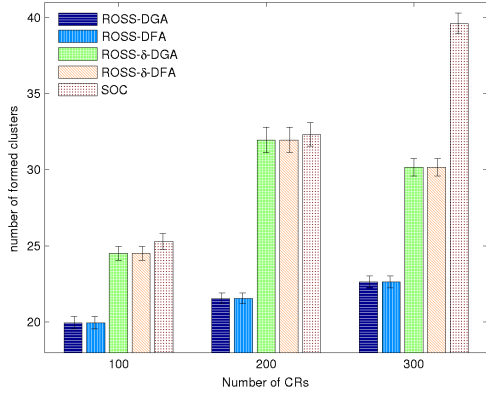


Fig. 19. The number of formed clusters.

Figures 20 21 22 illustrate the empirical cumulative distribution of CR nodes which reside in clusters with certain sizes in CRNs with different densities. When the variants of ROSS with size control feature are applied, the sizes of the most generated clusters are below δ , and most of them are around the 50% percentile. The sizes of clusters generated by ROSS-DGA and ROSS-DFA span a wider range than that with feature control feature. We find that average number of neighbours is roughly equal with the 95% percentile of the ROSS-DGA curve. As to SOC, the 95% percentiles are 36, 30, and 40. Overviewing the three Figures, we can see ROSS-DGA and ROSS-DFA show similar behaviour on cluster sizes. In contrary, the clusters generated from SOC demonstrate strong divergence on cluster sizes.

7 CONCLUSION

In this paper we design a distributed clustering scheme with the singleton congestion model, which forms robust clusters against

primary users' effect. Through simulation and theoretical analysis, we find that distributed scheme achieves similar performance with centralized optimization in terms of cluster survival ratio and number of control messages. This paper investigates the robust clustering problem in CRN extensively, and proves the NP hardness of this problem. A Light weighted clustering scheme ROSS is proposed, on the basis of which, we propose the fast version scheme and the scheme which generate clusters with desired sizes. These schemes outperform other distributed clustering scheme in terms of both cluster survival ratio and control overhead.

The shortage of ROSS scheme is it doesn't generate big clusters whose sizes exceed the cluster head's neighbourhood. This problem is attributed to fact that ROSS forms clusters on the basis of cluster head's neighbourhood, and doesn't involve interaction with the nodes outside its neighbourhood. In the other way around, forming big cluster which extends a cluster head's neighbourhood has limited applications, as multiple hop communication and coordination are required within such clusters.

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

Input: $d_j, g_j, j \in Nb_i \setminus CHs$. Empty sets τ_1, τ_2

Result: Returning 1 means i is cluster head, then d_j is set to 0, $j \in Nb_i \setminus CHs$. returning 0 means i is not CH.

```

1  if  $\nexists j \in Nb_i \setminus CHs$ , such that  $d_i \geq d_j$  then
2    | return 1;
3  end
4  if  $\exists j \in Nb_i \setminus CHs$ , such that  $d_i > d_j$  then
5    | return 0;
6  else
7    | if  $\nexists j \in Nb_i \setminus CHs$ , 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;
```

PROOF OF THEOREM 4.1

Proof. We consider a CRN which can be represented as a connected graph. (COMMENT: Why is it necessarily connected? No special case where the graph has two components?) To simplify the

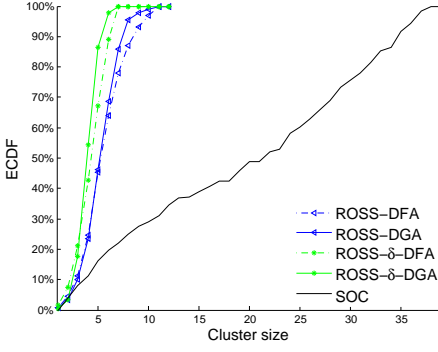


Fig. 20. 100 CRs, 30 PRs in network

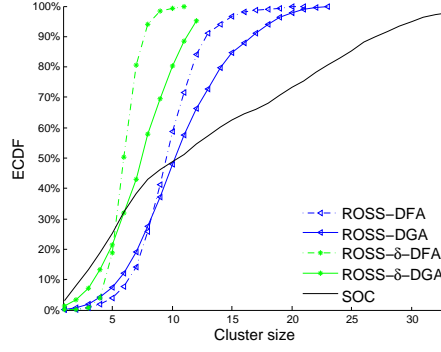


Fig. 21. 200 CRs, 30 PRs in network

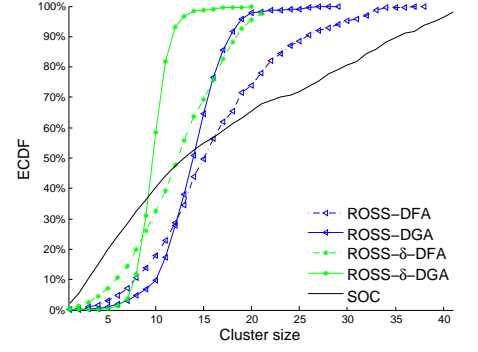


Fig. 22. 300 CRs, 30 PRs in network

Fig. 23. Cumulative distribution of CRs residing in clusters with different sizes

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

Input: Cluster C , empty sets τ_1, τ_2
Output: Cluster C has at least one CCC, or satisfies the requirement on cluster size

```

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

```

discussion, we assume the secondary users have unique individual connectivity degrees. Each user has an identical ID and a social connectivity degree. This assumption is fair as the social connectivity degrees and node ID are used to break ties in Algorithm 1, when the individual connectivity degrees are unique, it is not necessary to use the former two metrics. (COMMENT: What do you mean by breaking ties?? I don't see why your argument holds.)

For the sake of contradiction, let us assume there exist some secondary user α which is not included into any cluster. Then there is at least one node $\beta \in Nb_\alpha$ such that $d_\alpha > d_\beta$. Now, we distinguish between two cases: If β becomes cluster head, node α is included. If β is not a cluster head, i.e., β is not in any cluster (COMMENT: β not being a cluster head and β not belonging to any cluster is not the same statement!!!), we can repeat the previous analysis made on node α , and deduce that node β has at least one neighbouring node γ with $d_\gamma < d_\beta$. When no cluster head appears, this series (COMMENT: what is that? sequence?) of nodes with monotonically decreasing connectivity degrees might continue to grow. Finally, it ceases (COMMENT: not a good/suitable word!!)

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

Input: all claiming clusters $C \in S_i$

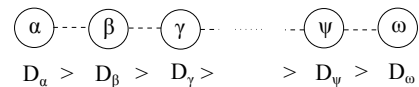
Output: one cluster $C \in S_i$, node i notifies all its claiming clusters in S_i about its affiliation decision.

```

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  $C \in S_i$ ,  $C = \arg \min(|K(C \setminus i)| - |K(C)|)$ 
3     then
4       return  $C$ ;
5   else
6      $\exists$  multiple  $C \in S_i$  which satisfies
7      $C = \arg \min(|K(C \setminus i)| - |K(C)|)$ ;
8      $\tau_1 \leftarrow C$ ;
9   end
10  if  $\exists$  only one  $C \in \tau_1$ ,  $C = \arg \max(K_{h_C} \cap K_i)$  then
11    return  $C$ ;
12  else
13     $\exists$  multiple  $C \in S_i$  which satisfies
14     $C = \arg \max(K_{h_C} \cap K_i)$ ;
15     $\tau_2 \leftarrow C$ ;
16  end
17  if  $\exists$  only one  $C \in \tau_2$ ,  $C = \arg \min |C|$  then
18    return  $C$ ;
19  else
20    return  $\arg \min_{C \in \tau_2} h_C$ ;
21 end

```

when either the individual connectivity degree is zero, or every secondary node has already been in the series (COMMENT: sequence?) of nodes. An example of the formed node series is shown as Figure 24.

Fig. 24. The node series discussed in the proof of Theorem 4.1, the deduction begins from node α

Now, we can see that the node ω is the last element of this

sequence. As ω does not have neighboring nodes with lower individual connectivity degree D , ω becomes the cluster head and incorporate all its one-hop neighbours, including the node before (COMMENT: you mean α ?) it in the series (here we assume that every newly formed cluster has common channels). After that, the node that joined a cluster will set its connection degree D to M , which enables the node further down in the list (WHAT ??? What is further down and which list??? I cannot follow.) to become a cluster head. In this way, all the nodes in the series are in at least one cluster in the inverse sequence (? inverse of?), which contradicts the assumption. Meanwhile, through this proof, we know that within at most N steps, all the nodes belong to certain clusters. (COMMENT: No, we don't! Where was that shown?) \square

PROOF OF THEOREM 5.1

Proof. We put the definition of weighted k -set packing problem here as it to be used in the analysis on the complexity of the centralised clustering problem.

DEFINITION 2: Weighted k -set packing.

Given a finite set $\mathcal{G} = \{g_1, \dots, g_N\} \subseteq \mathbb{N}$ of non-negative integers and a collection of sets $\mathcal{Q} = \{S_1, S_2, \dots, S_m\}$ such that $S_i \subseteq \mathcal{G}$ for every $1 \leq i \leq m$. Each set $S \in \mathcal{Q}$ is associated with a weight $w(S) \in \mathbb{R}$. Further, we are given a threshold value $\lambda \in \mathbb{N}$. The question is whether there exists a collection $\mathcal{S} \subseteq \mathcal{Q}$ such that \mathcal{S} contains only pairwise disjoint sets, i.e., for all $S, S' \in \mathcal{S}$ with $S \neq S'$ it holds that $S \cap S' = \emptyset$, and the total weight of the sets in \mathcal{S} is greater than λ , i.e., $\sum_{S \in \mathcal{S}} w(S) > \lambda$.

Weighted k -set packing is NP-hard when $k \geq 3$. [36]

To prove that the centralized clustering problem is NP-hard, we reduce the NP-hard problem *weighted k -set packing* to it to prove the former is at least as hard as the latter. We show the existence of a polynomial-time algorithm σ that transforms any instance \mathcal{S} (COMMENT: be careful - there is an overloading of \mathcal{S} of a weighted k -set packing into an instance $\sigma(\mathcal{S})$ of centralized clustering such that ... (to be completed - Erwin can do that)

W.I.o.g. let set $\mathcal{G} = \{1, \dots, N\}$.

ERWIN FINISHED HERE FOR THE MOMENT —

The polynomial algorithm σ consists of three transformations.

- Given a collection \mathcal{Q} , on basis of which we construct a CRN. We prepare N CR nodes who are labelled from 1 to N , put them on a 2 dimension space, and deem a pair of them can communicate if they appear in one same set in \mathcal{Q} . We regard each set in \mathcal{Q} is a cluster, whose number of CCCs equals to the weight of that set. Assuming two sets in \mathcal{Q} are $s_1 = \{1, 2\}$ and $s_2 = \{1, 2, 3\}$, then their weights are 3 and 5 respectively. We find it is impossible to map the sets into clusters in the same time, because the number of CCCs of the cluster which bases on s_1 should be no less than the cluster which bases on s_2 , as the latter has one extra node compared with the former cluster. But as to any instance of the solution to the weighted k -set packing problem, this contradiction doesn't happen because the instance \mathcal{S} contains only disjoint sets, thus at most only one set of s_1 and s_2 appears in \mathcal{S} .
- In the second step, we transform the instance \mathcal{S} to \mathcal{S}' by adding dummy elements into each set in \mathcal{S} . For each set $s_i \in \mathcal{S}$, the elements in s_i are duplicated, for instance, given $s_i = \{1, 4, 6\}$, the dummy set $s'_i = \{1, 1, 4, 4, 6, 6\}$. The purpose of this transformation is to eliminate the set in \mathcal{S} , which has single element. The weight of set is unchanged after this transformation, i.e., $\omega(s_i) = \omega(s'_i)$. This transformation requires $\sum_{s_i \in \mathcal{S}} |s_i|$ steps.

- In this step, we transform the instance \mathcal{S}' to a clustering solution for CRN. We prepare a second pool of CR nodes which are identical with the CR nodes prepared in step 1, i.e., identical IDs and channel availabilities on them, we call these CR nodes as dummy nodes. We locate these CR nodes besides the CR nodes with the same IDs in the CRN built in step 1, and there is connection between the CR node and its dummy node (the one CR node and its dummy node can be seen as two transceivers at one node). Because of the dummy nodes, the clustering solution which corresponds to \mathcal{S}' doesn't have singleton cluster. This transformation requires $2 \cdot \sum_{s_i \in \mathcal{S}} |s_i|$ steps. Afterwards, the CR node whose ID doesn't appear in any set in \mathcal{S} becomes single node clusters, according to the definition of clustering problem in CRN, the number of CCCs in these single node clusters is 0. These singleton clusters and the clusters in \mathcal{S} constitute a clustering solution, and finding the singleton clusters requires at most N steps. An example is shown in Table 4.

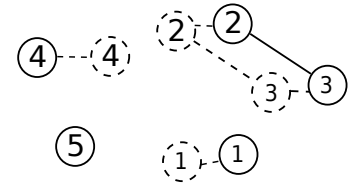
| | |
|--|---|
| \mathcal{N} | $\{1, 2, 3, 4, 5\}$ |
| \mathcal{Q} | $\{(1), (1, 5), (1, 2, 4), (2, 3), (4)\}$ |
| Instance for Weighted k -set packing | $\{(1), (2, 3), (4)\}$ |
| Instance with dummy elements | $\{(1, 1), (2, 2, 3, 3), (4, 4)\}$ |
|  | |

TABLE 4

We have crossed the hurdle of finding one polynomial algorithm σ which transforms an instance of weighted k -set packing to an instance of the clustering problem in CRN. Now we look into the step 2 in the reduction. As to an instance \mathcal{S} for weighted k -set packing, the sum weights is identical to the sum of CCCs in the CRN mapped from \mathcal{S}' , even \mathcal{S} contains set which only has one element. Thus, when the instance \mathcal{S} is one solution and its sum weights is greater than λ , in the CRN which is mapped from \mathcal{S}' , the summed number of CCCs of the clusters is greater than λ . When there is no solution out of set \mathcal{G} for weighted k -set packing, the summed number of CCCs of the clusters in the mapped CRN is also smaller than λ .

Thus, weighted k -set packing can be reduced to clustering problem in CRN, then the latter problem is of NP-hard. \square

ACKNOWLEDGMENT

The authors would like to thank xxxx

REFERENCES

- [1] I. Mitola, J. and J. Maguire, G.Q., "Cognitive radio: making software radios more personal," *Personal Communications, IEEE*, vol. 6, no. 4, pp. 13–18, aug. 1999.
- [2] Q. Zhao and B. Sadler, "A survey of dynamic spectrum access," *Signal Processing Magazine, IEEE*, vol. 24, no. 3, pp. 79–89, May 2007.

- [3] A. Sahai, R. Tandra, S. M. Mishra, and N. Hoven, "Fundamental design tradeoffs in cognitive radio systems," in *Proc. of the First International Workshop on Technology and Policy for Accessing Spectrum*, ser. ACM TAPAS '06.
- [4] I. F. Akyildiz, B. F. Lo, and R. Balakrishnan, "Cooperative spectrum sensing in cognitive radio networks: A survey," *Phys. Commun.*, vol. 4, no. 1, pp. 40–62, Mar. 2011.
- [5] C. Sun, W. Zhang, and K. B. Letaief, "Cluster-based cooperative spectrum sensing in cognitive radio systems," in *proc. of IEEE ICC*, 2007, pp. 2511–2515.
- [6] J. Zhao, H. Zheng, and G.-H. Yang, "Spectrum sharing through distributed coordination in dynamic spectrum access networks," *Wireless Com. and Mobile Computing*, vol. 7, no. 9, pp. 1061–1075, 2007.
- [7] D. Willkomm, M. Bohge, D. Hollós, J. Gross, and A. Wolisz, "Double hopping: A new approach for dynamic frequency hopping in cognitive radio networks," in *Proc. of PIMRC 2008*.
- [8] C. Passiatore and P. Camarda, "A centralized inter-network resource sharing (CIRS) scheme in IEEE 802.22 cognitive networks," in *Proc. of IFIP Annual Mediterranean Ad Hoc Networking Workshop 2011*, 2011.
- [9] A. A. Abbasi and M. Younis, "A survey on clustering algorithms for wireless sensor networks," *Comput. Commun.*, vol. 30, no. 14–15, pp. 2826–2841, 2007.
- [10] H. D. R. Y. Huazi Zhang, Zhaoyang Zhang¹ and X. Chen, "Distributed spectrum-aware clustering in cognitive radio sensor networks," in *Proc. of GLOBECOM 2011*.
- [11] B. E. Ali Jorio, Sanaa El Fkihi and D. Aboutajdine, "An energy-efficient clustering routing algorithm based on geographic position and residual energy for wireless sensor network," *Journal of Computer Networks and Communications*, vol. 2015, April 2015.
- [12] Q. Wu, G. Ding, J. Wang, X. Li, and Y. Huang, "Consensus-based decentralized clustering for cooperative spectrum sensing in cognitive radio networks," *Chinese Science Bulletin*, vol. 57, no. 28–29, pp. 3677–3683, 2012.
- [13] V. Kawadia and P. R. Kumar, "Power control and clustering in ad hoc networks," in *Proc. of INFOCOM '03*, 2003, pp. 459–469.
- [14] C. R. Lin and M. Gerla, "Adaptive clustering for mobile wireless networks," *IEEE Journal on Selected Areas in Communications*, vol. 15, pp. 1265–1275, 1997.
- [15] S. Basagni, "Distributed clustering for ad hoc networks," *Proc. of I-SPAN '99*, pp. 310–315, 1999.
- [16] J. Sućec and I. Marsic, "Clustering overhead for hierarchical routing in mobile ad hoc networks," in *Proc. of IEEE INFOCOM 2002*, vol. 3, 2002, pp. 1698 – 1706 vol.3.
- [17] H. Wu and A. Abouzeid, "Cluster-based routing overhead in networks with unreliable nodes," in *Proceedings of IEEE WCNC 2004*, vol. 4, march 2004, pp. 2557 – 2562 Vol.4.
- [18] T. Chen, H. Zhang, G. Maggio, and I. Chlamtac, "Cogmesh: A cluster-based cognitive radio network," *Proc. of DySPAN '07*.
- [19] D. Wu, Y. Cai, L. Zhou, and J. Wang, "A cooperative communication scheme based on coalition formation game in clustered wireless sensor networks," *IEEE Transactions on Wireless Communications*, vol. 11, no. 3, pp. 1190–1200, march 2012.
- [20] K. Baddour, O. Ureten, and T. Willink, "Efficient clustering of cognitive radio networks using affinity propagation," in *Proc. of ICCCN 2009*.
- [21] A. Asterjadhi, N. Baldo, and M. Zorzi, "A cluster formation protocol for cognitive radio ad hoc networks," in *Proc. of European Wireless Conference 2010*, pp. 955–961.
- [22] L. Lazos, S. Liu, and M. Krunz, "Spectrum opportunity-based control channel assignment in cognitive radio networks," *Proc. of SECON '09*, pp. 1–9, jun. 2009.
- [23] S. Liu, L. Lazos, and M. Krunz, "Cluster-based control channel allocation in opportunistic cognitive radio networks," *IEEE Trans. Mob. Comput.*, vol. 11, no. 10, pp. 1436–1449, 2012.
- [24] B. Clark, C. Colbourn, and D. Johnson, "Unit disk graphs," *Annals of Discrete Mathematics*, vol. 48, no. C, pp. 165–177, 1991.
- [25] Y. Zhang, G. Yu, Q. Li, H. Wang, X. Zhu, and B. Wang, "Channel-hopping-based communication rendezvous in cognitive radio networks," *IEEE/ACM Transactions on Networking*, vol. 22, no. 3, pp. 889–902, June 2014.
- [26] Z. Gu, Q.-S. Hua, and W. Dai, "Fully distributed algorithms for blind rendezvous in cognitive radio networks," in *Proceedings of the 2014 ACM MobiHoc*, ser. MobiHoc '14.
- [27] Y. P. Chen, A. L. Liestman, and J. Liu, "Clustering algorithms for ad hoc wireless networks," in *Ad Hoc and Sensor Networks*. Nova Science Publishers, 2004.
- [28] E. Perevalov, R. S. Blum, and D. Safi, "Capacity of clustered ad hoc networks: how large is "large"?" *IEEE Transactions on Communications*, vol. 54, no. 9, pp. 1672–1681, Sept 2006.
- [29] A. MacKenzie and S. Wicker, "Game theory in communications: motivation, explanation, and application to power control," in *Proc. of IEEE GLOBECOM 2001*.
- [30] J. O. Neel, "Analysis and design of cognitive radio networks and distributed radio resource management algorithms," Ph.D. dissertation, Blacksburg, VA, USA, 2006, aA13249450.
- [31] B. Wang, Y. Wu, and K. R. Liu, "Game theory for cognitive radio networks: An overview," *Comput. Netw.*, vol. 54, no. 14, pp. 2537–2561, Oct. 2010.
- [32] B. J. S. Chee and C. Franklin, Jr., *Cloud Computing: Technologies and Strategies of the Ubiquitous Data Center*, 1st ed. Boca Raton, FL, USA: CRC Press, Inc., 2010.
- [33] H. Ackermann, H. R  uglin, and B. V  ucking, "Pure Nash equilibria in player-specific and weighted congestion games," *Theoretical Computer Science*, vol. Vol. 410, no. 17, pp. 1552 – 1563, 2009.
- [34] X.-Y. Li, Y. Wang, and Y. Wang, "Complexity of data collection, aggregation, and selection for wireless sensor networks," *IEEE Transactions on Computers*, vol. 60, no. 3, pp. 386–399, 2011.
- [35] M. Onus, A. Richa, K. Kothapalli, and C. Scheideler, "Efficient broadcasting and gathering in wireless ad-hoc networks," in *Proc. of ISPAN 2005*.
- [36] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.

PLACE
PHOTO
HERE

Di Li received BE and MS degrees in control engineering from Zhejiang University and Shaanxi University of Science and Technology respectively in China. He worked with James Gross for his PhD in RWTH Aachen University since 2010.

PLACE
PHOTO
HERE

Erwin Fang Biography text here. Biography text here. Biography text here. Biography text here. Biography text here.

PLACE
PHOTO
HERE

James Gross Biography text here. Biography text here. Biography text here. Biography text here. Biography text here.