

# CCIT Report #780 January 2011

# Cross-Entropy Optimized Cognitive Radio Policies

Boris Oklander Department of Electrical Engineering Technion – Israel Institute of Technology Haifa, Israel oklander@tx.technion.ac.il Moshe Sidi Department of Electrical Engineering Technion – Israel Institute of Technology Haifa, Israel moshe@ee.technion.ac.il

Abstract—In this paper we consider cognitive processes and their impact on the performance of cognitive radio networks (CRN). We model the cognition cycle, the main control process of the cognitive radio (CR), during which CR sequentially senses and estimates the environment state, creates plans based on the knowledge (models) of itself and that of the environment, makes decisions in order to optimize certain objectives and then acts. The proposed framework is analyzed and the performance of the CRN is evaluated. We show the impact of the sensing rate and the system dynamics on the waiting times of secondary users. Then model-based analysis is used to solve control and decision making tasks, which actually gives the radio its "cognitive" ability. Particularly, we design an efficient strategy for accessing the vacant spectrum bands and managing the transmission-sampling trade-off. In order to cope with the high complexity of this problem the policy search uses the stochastic optimization method of cross-entropy. The developed cognition cycle model represents CRN ability to intelligently react to external environment and internal state changes and gives a good understanding of the cross-entropy optimized policies.

Keywords-cognitive radio networks; dynamic spectrum access; state estimation; queueing analysis; cross-entropy

# I. INTRODUCTION

Cognitive Radio Networks (CRN) provide new horizons to the next generation of wireless communications. This new communication paradigm is a candidate to cope with a wide spectrum of challenges arising in the face of the increasing demand for wireless access in voice, video, multi-media and other high rate data applications. Although researchers and standardization bodies generally agree that CR should be able to sense the environment and autonomously adapt to changing operating conditions, there are different views concerning the levels of cognitive functionality [1]. CRNs are envisioned to aid both the user and the network to mitigate the growing communication demands by their advanced capabilities represented by the cognition cycle [2].

Cognition cycle is the main control process which enables CR to stay aware of its communication environment and to adapt to its changing conditions. There are different views of what phases the cognition cycle consists [2],[4], but basically all the versions share the phases of Boyd's observation, orientation, decision and action (OODA) loop [5] (see Fig. 1). During the observation phase, CR continuously senses the environment in order to collect the input information for the cognition cycle. In the orientation phase, CR uses the gathered information from its sensors along with its general knowledge (models) of the communication environment to estimate the current network state. Next, given the estimated network conditions, CR enters the decision-making phase in which it applies some policy to decide on the course of action. The CR's policy comprises the knowledge of both the environment and the CR, and it is optimized to meet communication goals. Finally, CR completes the cognition cycle by entering the action phase, which carries out the chosen actions. In addition, machine learning can be structured into these phases of cognition cycle in order to update them in the face of changing environment's situation or user needs.



Figure 1. Cognition Cycle: Observe, Orient, Decide and Act phases.

The cognition cycle implies strong correlation between the perception (sensing and estimation) and the action (transmissions). Essentially, the interdependence of perception and action is a fundamental principle governing CRs behavior. Perceiving both the communication environment and the self states enables CR to intelligently adapt its actions in the face of the dynamically changing conditions of the network. CR has some control over the sensing process, and therefore it can deliberately modify its perception level by changing the resource allocation (e.g. varying the sampling rate). Both the perception and the action processes make use of CRs limited resources such as computation power, spectrum bandwidth etc. Therefore, by applying appropriate policy for resources allocation, CR should adaptively optimize its operating point.

Different studies have addressed CRNs capability of opportunistic spectrum access [3],[4],[6], in which spectrum bands licensed to primary users (PU) are shared with the cognitive users called secondary users (SU). It is well known that a significant part of the allocated spectrum is vastly underutilized [7],[8],[9], and the CRN goal in this scheme is

to improve spectrum utilization while avoiding interference with the PUs [10]-[12]. This requires management of the sampling-transmission tradeoff [13]-[16].

It is not common to find studies that directly address the interdependent processes composing the cognition cycle. The main reason for this is the difficulty to design analytically tractable models for systems characterized by cognitive behavior. As for now, a substantial gap remains between the perception and action-taking models. In [17] state of the art protocols for medium access in cognitive radio networks are overviewed. The authors point out that the existing works do not fully integrate both the spectrum sensing and the spectrum access in one framework which is required in order to maintain the capability of adaptation to the environment changes [18].

The authors of [19] derive a threshold strategy for the sequential channel sensing process aiming to maximize the aggregated throughput of CRN. While the model in [19] assumes independent transmissions over different channels, our model can deliberately utilize any number of channels it can observe simultaneously and therefore achieves higher degree of spectral agility at the expense of strong correlation between the channels. Another important difference is that our model embeds the CRs buffer, which allows a more accurate performance evaluation of CR in general and obtaining the delay performance in particular. The authors of [20] derive a queueing framework to study the performance of CRN accessing the spectrum in an opportunistic manner. Although this model allows an analytic study of CRN performance, it lacks the modeling of the cognition cycle as it neglects the phase of environment sensing and its state estimation. The model in [21], lacks the sampling-transmission tradeoff and penalty for interference with PU. Our analytic model inherently combines these important processes.

This paper contributes a unified model for the cognition cycle. We model the environment as a stochastic system dynamically changing its state. CRN observes the environment stochastically in order to estimate the state of PU network. The perception process is an integral part of the system's overall behavior, and we solve simultaneously for the optimal control of the observations and the spectrum access. In this way, CRN is able to allocate the spectrum resources according to the overall task requirements. We model the CRN buffer as an infinite queue and seamlessly embed it in our analytic framework. A few factors determine the instantaneous service rate of the CRN queue. These factors include the number of accessed channels, proportion of the bandwidth assigned for transmission (residual part goes to sampling) and penalty for accessing channels occupied by PU.

An additional contribution of this paper is the introduction of cross-entropy optimized policy for controlling the CRN. Given a certain control policy, the described above framework is solved using matrix-geometric approach [21]. However, the task of policy optimization is rather hard due to the high complexity of the model. To overcome this problem we use the method of stochastic optimization of cross-entropy [26], which is an efficient tool at hand for the task of policy optimization [26][27],[28]. The resulting policies reflect the intelligent behavior induced by the described above cognition cycle.

This paper is organized as follows. In section II we present a stochastic model of the cognition cycle. In section III the model is analyzed using Matrix-Geometric approach, which is applicable to vector-state Markov processes that have repetitive structure. The numerical analysis of the proposed model introduces insights into the performance of

CRN. In particular, we establish the relations between different quantities such as input rate, environment state estimation rate and waiting times of SU. Section IV we use the cross-entropy method to optimize the control policy responsible for allocation of the CRN resources. Section V summarizes the work.

# II. COGNITION CYCLE MODEL

We regard the cognition cycle as an aggregation of interdependent processes through which CRN interacts with the communication environment. CRN access channels temporarily unoccupied by PU in order to transmit data. We denote by  $S_t$  the state of the environment at time t, which is actually the number of available channels for CRN access. In the case when CRN tries to access channels erroneously estimated as vacant, the transmissions fail. This penalty for interference with PU implies a significant incentive for CRN to allocate resources required for enhancing its perception level.

The perception process consists of sensing the environment and estimating its state. We denote by  $\hat{S}_t$  the estimation of the environment state  $S_t$ . CRN observes the environment by sampling the network channels and it has some control over the observation process by deliberately tuning the sampling rate  $\delta_t$  over time. For example, CRN could increase the sample rate in order to keep track of rapidly changing network states characterized by high throughput potential while decreasing it for slowly changing states. Since CRN estimates the network state, we assume that the sampling rate may depend on  $\hat{S}_t$ . It is reasonable to assume that due to the physical and the hardware limitations, the transmissions and the observations are mutually exclusive and hence the sampling and the transmission rates are negatively correlated. This condition forms the throughput-sampling tradeoff which was the focus of the study in [13]-[16] and is an intrinsic part of our model as we point out later.

In the following subsections, we model the environment's dynamics, the cognition cycle and the CRN data transmission process. Then, we unify these models under the entire system framework. Closing the loop makes it possible to analyze the cognition cycle and to evaluate the performance of the CRN. We assume that CRN knows the correct models of both the environment and the transmitter. This assumption reduces the need of updating the models (for example through machine learning methods) and allows us to focus on modeling the perception and decision making phases of the cognition cycle.

#### A. Environment Model

In the scenario under consideration, CRN accesses the network channels in an opportunistic manner to create virtual unlicensed bands, i.e., bands that are shared with PU on a non-interfering basis. We consider a general scenario of wireless communication system which consists of M channels. There are M PUs in the system, while every PU has an exclusive access to a single channel. Every PU alternates between transmitting and idle states. The ON (OFF) period of a channel corresponds to the time interval  $T_{ON}$  ( $T_{OFF}$ ) during which a PU is transmitting (idle). We assume that  $T_{ON}$  and  $T_{OFF}$  intervals are exponentially distributed with parameters  $\alpha$  and  $\beta$ , respectively.

CRN uses the channels to form a pool of M spectral bands. In this mode of operation CRN look for "holes" in the spectrum and dynamically adapt its transmissions over

unused bands. Note that the holes do not have to be contiguous [23]. Additionally, once CRN detects PU appears in a frequency band all SU leave this band immediately, giving priority to PU and avoiding interference. Since the PU are statistically independent, the number of bands available for SU access  $S_t$  ( $S_t \in \{0,1,...,M\}$ ) at time t is a birth-death process with birth-rate (M-m) $\alpha$  and death-rate  $m\beta$  when  $S_t=m$ ,  $m \in \{0,1,...,M\}$  (see Fig. 2).



Figure 2. Aggregate birth-death process of unoccupied bands.

#### B. Perception Model

Here we model the perception process, which is an aggregation of the observation and the orientation phases of the cognition cycle. The environment state  $S_t$  is unknown and therefore CRN estimates it through sensing. As was assumed before, CRN knows the environment model and its parameters. In our case of structured environment model the parameters are  $\alpha$ ,  $\beta$  and M. CRN uses the environment model and the data from sensors to obtain the estimation  $\hat{S}_t$ , which is the output of the unified perception phase of the cognition cycle. The perception process updates  $\hat{S}_t$  at random time instants  $t_k$ ,  $k \in \{0,1,2...\}$ . We assume that the time it takes to update the estimation is exponentially distributed. CRN adaptively tunes the update rate  $\delta_t$  according to its current estimate  $\hat{S}_t$ . The notations of  $\hat{S}_k$  and  $S_k$  describe the values of  $\hat{S}_t$  and  $S_t$  at time  $t_k$ . At each instant  $t_k$ , the estimation  $\hat{S}_k$  is updated to be the true value of  $S_k$  and remains unchanged till the next update instant  $t_{k+1}$ .



Figure 3. CTMC of the  $Z_t = \{\hat{S}_t, S_t\}$  process for M=3. The horizontal transitions describe the changes of network state  $S_t$ . The vertical transitions describe the updates of the estimator  $\hat{S}_t$  to the correct value of  $S_t$ .

The compound process  $Z_t = {\hat{S}_t, S_t}$  describes the mutual evolvement of both the environment and the estimation processes which can be shown to be a continuous time Markov chain (CTMC) (see Fig. 3). In this CTMC the horizontal transitions describe the changes of environment state  $S_t$ . The vertical transitions describe the updates of the estimator  $\hat{S}_t$  toward the correct value of  $S_t$ . Note that the states for which  $S_t = \hat{S}_t$  act as absorbers of the vertical transitions. Once the process enters such a state, the vertical transitions hold off until the moment when the environment state changes.

#### C. Decision Making

The decision-making phase of the cognition cycle employs some policy P for both transmission-sampling tradeoff management and for channels allocation. As we already mentioned, at any instant CRN either senses or transmits over a channel. The transmission rate of SU over a single unoccupied by PU channel is  $\mu$  [bit/sec]. We introduce the tradeoff parameter  $\theta$  ( $0 \le \theta \le 1$ ) which divides the available bandwidth between the transmissions and sampling, where the portion  $\theta$  of the channel is assigned for transmission and the remaining part  $(1-\theta)$  is assigned for the sampling process. For a given value of  $\theta$ , the effective transmission rate over a single channel is therefore  $\theta\mu$  [bit/sec] and the resulting update rate of the estimations is  $(1-\theta)\mu B$  [1/sec]. The constant 1/B [bit] is the number of bits required for updating the estimation  $\hat{S}_t$  and it is subject to the physical layer issues. The other responsibility of the policy P is the channel allocation  $C_t$ , which is the number of channels over which CR tries to transmit at time t.

In this work, we consider state-dependent policies, meaning that the decisions are made based on the estimation of the network state  $\hat{S}_t$  and the internal buffer state  $X_t$ . The internal buffer state  $X_t$  is the number of SU packets waiting for transmission at time t. For the sake of simplicity, in the following modeling we assume that CRN makes decisions based on a greedy policy  $P_G$ ,  $C_t=P_G(\hat{S}_t,X_t)$ . The greedy policy aims to increase the throughput by scheduling transmissions over all the channels that are estimated as unoccupied by PU while keeping constant tradeoff parameter:

$$C_{t} = P_{G}(\hat{S}_{t}, X_{t}) = \begin{cases} \hat{S}_{t} & X_{t} > 0\\ 0 & X_{t} = 0 \end{cases}$$
(1)

This assumption of greedy policy is removed later in section IV when we search for optimized policies in order to achieve better CRN performance.

#### D. Transmission Process

The arrivals generated by SU are modeled as a Poisson process with rate  $\lambda$  [bit/sec] and service time exponentially distributed with rate  $\mu_t$  [bit/sec], which changes with time dependent on a few factors. These factors are the number of accessed channels  $C_t$ , the proportion of the channels bandwidth allocated for transmission  $\theta$ , the actual state of the environment  $S_t$  and the penalty for interfering with PU. The combination of these factors results in

$$\mu_t = \begin{cases} \theta C_t \mu & C_t \leq S_t \\ 0 & C_t > S_t \end{cases}$$
(2)

It can be seen from (1) that when the decisions are made according to the greedy policy  $P_G$ , we may substitute  $\hat{S}_t$  for  $C_t$  since transmissions occur only for  $X_t>0$ . As it can be seen from (2), our model introduces penalty for CRN when it accesses channels that are in use of PU ( $C_t>S_t$ ). This type of service models the opportunistic spectrum access of CRN giving the highest priority to PUs. For example, during the periods when all the bands are occupied by PUs ( $S_t=0$ ) no CRN packets are transmitted independently of  $\hat{S}_t$ .



Figure 4. Illustration of the CTMC of the CRN model. The transitions in the  $(\hat{S}_t, S_t)$  plane are identical to those in Fig. 3. The transitions between the levels of the process (along the  $X_t$  axis) are ommited here in sake of keeping visuability, they are presented in the Appendix A.

#### E. System Process

Now we aggregate the environment dynamics, the cognition cycle and the transmission process into a unified system model. We define  $\{X_t, Z_t\}$  to be the process of the entire system for which at time *t* there are  $X_t$  ( $X_t \in \{0, 1, 2, ...\}$ ) queued packets of SU, which is the level of the process, and  $Z_t = \{\hat{S}_t, S_t\}$  ( $Z_t \in \{0, ..., M\} \times \{0, ..., M\}$ ), which is state within the level. This process forms a three dimensional CTMC illustrated in Fig. 4, which is homogeneous, irreducible and stationary.

The exact structure of transitions within the CTMC and its analysis by means of matrix geometric approach are presented in Appendix A. In the next section, we use the results of the analysis of the CTMC to evaluate the performance of the CRN described by our model.

## **III. CRN PERFORMANCE EVALUATION**

We first aim at evaluating the performance of the estimator  $\hat{S}_t$ . The mean square error (MSE) of an estimator is a common way to evaluate its performances. MSE quantifies the difference between an estimator and the true value of the quantity being estimated. In our case

$$MSE = E\left[\left(\hat{S}_{t} - S_{t}\right)^{2}\right] = \sum_{i,j=0}^{M} (i - j)^{2} \Pr\left(\hat{S}_{t} = i, S_{t} = j\right)$$
(3)

The probabilities  $Pr(\hat{S}_t=i, S_t=j)$  can be easily obtained by solving a CTMC of the  $Z_t$  process, like the one presented in Fig 3.



Figure 5. MSE of  $\hat{S}_t$  for parameters M=5,  $\alpha=0.5$ ,  $\beta=1$ ,  $\mu=1$ ,  $\lambda=1$ ,  $k=\{10,100,1000,10000\}$ . The MSE improves for growing values of  $\delta_t$ .

An interesting observation, which characterizes the performance of CRN, is its dependence not only on the fraction of time available for SU to access the channel, but also on the pattern of spectrum usage of PU. Let  $\gamma_k = (k\alpha)/(k\beta)$ , k>0. It is obvious that  $\gamma_k = \gamma$  and therefore the fraction of time that the channel is available for the SU remains constant for all k>0. The difference k causes is in the pattern of spectrum usage of PU. For low values of k the rates are slow and PU are characterized by a persistent behavior in which they remain in transmitting or idle states for long periods of time compared to SU. When k is high, PU behave in an oscillatory manner alternating quickly between the transmitting and idle states. In Fig. 5 it can be seen that MSE of the estimator  $\hat{S}_t$  improves for increasing values of  $\delta_t$ . As the PU oscillate more frequently (increased k) the update rate  $\delta$  should be significantly increased in order to keep the same MSE value.



Figure 6. Waiting time vs. k for different values of  $\gamma$ . Parameter values  $\alpha=1$ ,  $\beta=2$ ,  $\mu=1$ ,  $\lambda$  varies in the stable region of the system  $\rho=\lambda/\mu \leq M\alpha/(\alpha+\beta)$ .

Next we assess the communication performance of the SU. We focus on the waiting time of SU. In Appendix A a way to obtain  $\pi_0$  and *R* is presented. Using these quantities it is possible to calculate the number of queued packets of SU, denoted by  $N_q$ :

$$N_{q} = \sum_{j=1}^{\infty} j\pi_{j} e = \pi_{1} \sum_{j=1}^{\infty} jR^{j-1} e = \pi_{0} R(I-R)^{-2} e$$
(4)

where  $\pi_i = (\pi_{i,1}, \pi_{i,2}, ..., \pi_{i,M^2})$  and  $\pi_{i,j}$  are the stationary probabilities of the process  $\{X_t, Z_t\}$  to be at level *i* and state *j* within that level. Using  $N_q$  and Little's law, we can obtain the waiting time of SU:

$$W = N_q / \lambda = \left( \pi_0 R (I - R)^{-2} e \right) / \lambda$$
(5)



Figure 7. Waiting times of CTMC model for parameters M=5,  $\alpha=0.5$ ,  $\beta=1$ ,  $\mu=1$ ,  $\lambda=1$ ,  $\delta=\{10,100,1000,10000\}$ .

We examine the case when the estimation is perfect, i.e.  $\hat{S}_t=S_t$  for all *t*. This situation is achieved for  $\delta_t \rightarrow \infty$ ,  $\forall \hat{S}_t$ . The resulting performance of SU dependent on *k* is presented in Fig. 6. It is noticeable that for persistent behavior of the PU, the CRN performance weaken and SU have to wait longer periods on average although the channel is available for secondary access the same fraction of time.

Next we examine the behavior of the waiting time for finite update rates, see Fig. 7. When the update rate  $\delta_t (\forall \hat{S}_t)$  is significantly higher than the transition rates  $(\alpha, \beta)$ , the estimator  $\hat{S}_t$  is characterized by a small MSE (see Fig. 8). As a result the curves coincide in the corresponding interval of k values. In this case the waiting time behaves in the same manner as if the estimation process had small MSE. However, when the transition rates grow, the performance of CRN becomes sensitive to the estimation process. As it can be seen from Fig. 7, for k>1, the accuracy of the estimation process affects the performance of SU significantly. When the updates of  $\hat{S}_t$  occur too slowly compared to the environment dynamics, the waiting time increases. Each of the curves describes longer waiting times dependent on  $\delta$ . Further, it can be seen from the graphs that the waiting time saturates in a rapidly changing environment, however the system remains stable. This can be explained by the fact that when the environment state fluctuates quickly, the probability  $\Pr(\hat{S}_t=i,S_t=j)$  remains positive and independent of k or  $\delta$ , which can be seen from Fig 8.



Figure 8. MSE of  $\hat{S}_t$  for parameters *M*=5,  $\alpha$ =0.5,  $\beta$ =1,  $\mu$ =1,  $\lambda$ =1,  $\delta$ ={10,100,1000,10000}.

As a summary of the analysis we plot in Fig. 9 the performance curves of CRN for different values of M assuming perfect estimation. It is clear that the performance for different systems (different M values) saturate when  $\lambda$  approaches high values. It is interesting to notice that from this plot one can learn about system trade-offs. For example, one can answer the question whether better performance could be reached by splitting the SU in two groups generating half the original traffic rate (0.5 $\lambda$ ) and with separated spectrum pools of M/2 channels. Comparing the waiting times at the point  $\lambda$ =6 on curve M=8 to the point  $\lambda$ =3 on curve M=4 shows that using a larger spectrum pool improves the performance.



Figure 9. Waiting times of CTMC model for  $M=\{4,6,8,10\}, \alpha=10, \beta=2, \mu=1$ .

# IV. CRN POLICY OPTIMIZATION

In previous sections, we modeled the cognition cycle, analyzed it and evaluated its impact on the performance of the CRN. The evaluation was carried out for CRN that makes decisions based on some arbitrarily chosen greedy policy  $P_G$ . In this section, we aim to improve the performance of the cognition cycle and the CRN by optimizing the decision making process, i.e., by optimizing the policy.

# A. Problem formulation

In our framework, a policy *P* governs the decision-making phase of the cognition cycle. This policy is responsible for managing the sampling-transmission tradeoff by tuning the continuous parameter  $\theta_t$ , and for allocation of channels,  $C_t$ . The values of  $C_t$  and  $\theta_t$  are determined dependently on current estimation of the networks state  $\hat{S}_t$ , current CRN buffer state  $X_t$ , entire system model and its parameters, which we denote by  $\Omega$ :

$$(C_t, \theta_t) = P(S_t, X_t; \Omega) \tag{6}$$

We aim to optimize CRN performance by minimizing the average waiting time *W* of SU. In the previous section, we calculated *W* by applying the matrix geometric analysis to the 3-D CTMC and the Little's law. The 3-D CTMC structure embeds the policy *P* as follows: the levels transitions ( $X_t$ ) are affected by the service rate  $\mu_t$  (eq. 2) and the state transitions ( $Z_t$ ) within the level are affected by the estimation update rate  $\delta_t$  given by  $\delta_t = (1-\theta_t)\mu B$ . Therefore, given the system structure and its parameters  $\Omega$ , we regard the average waiting time *W* of SU as a function of the policy *P*,  $W=W(P;\Omega)$ . The resulting optimization problem is given by: where  $\Pi$  is the set of all the feasible policies, i.e., policies which for valid inputs  $\hat{S}_t \in \{0, ..., M\}$  and  $X_t \in \{0, 1, 2, ...\}$  decide on valid values for the of  $\theta \in [0, 1]$  and  $C_t \in \{0, 1, 2, ..., M\}$ . Our optimization problem (7) is a complicated one. First, it can be shown that the problem is not convex, and the gradient-based techniques are not applicable since it is difficult to obtain a gradient for W. Next, the set  $\Pi$  consists of policies comprising both continuous ( $\theta$ ) and discrete ( $C_t$ ) action spaces, which requires special approach for optimization. Additionally, the problem exhibits a high computational complexity, due to the rapidly growing (with M) set of feasible policies  $\Pi$ .

We solve this problem by applying the cross-entropy (CE) method of stochastic optimization. CE method is a state-of-the-art method for solving combinatorial and multi-extremal optimization problems. In the following subsection, we review briefly the CE method and demonstrate its application for our optimization problem. The readers interested in further details are referred to [26].

### B. Cross-Entropy based Stochastic Optimization

The main idea behind the CE method is to define for the original optimization problem an associated stochastic problem and then to solve efficiently the associated problem by an adaptive scheme. The described below procedure sequentially generates random solutions which converge stochastically to the optimal or near-optimal one.

We define a stochastic policy  $P((C_t, \theta_t) | \sigma(\hat{S}_t, X_t))$  as the associated stochastic problem for (7).  $P((C_t, \theta_t) | \sigma(\hat{S}_t, X_t))$  is the probability of choosing action  $(C_t, \theta_t)$  when CRN's state is  $(\hat{S}_t, X_t)$  according to the parameter  $\sigma(\hat{S}_t, X_t)$ . In the following we use shorthand notation of  $\sigma$  for  $\sigma(\hat{S}_t, X_t)$ . For the defined associated stochastic problem, the CE method iteratively draws sample policies  $P^{(k)}$  (k=1,2,...,K) from the defined above probability and calculates the average waiting time  $W(P^{(k)}; \Omega)$  for each sample. Then N (N < K) best samples graded by their related average waiting time are used to update the parameters  $\sigma$ , in order to produce better samples in the next iteration. The algorithm stops when the score of the worst selected sample no longer improves significantly. The exact CE algorithm is presented in Appendix B.

# C. Cross-Entropy Optimized Policies

We present here policies obtained from CE optimization and examine them in order to get insights concerning the optimal decision-making process in CRN. As in the previous sections we are interested to reveal the impact of the cognition cycle and the dynamics of the environment on the optimal policy. We set the parameters of the environment ( $\Omega$ ): the number of PU channels is M=6, and the transmission rate over every channel is  $\mu=1$ , the constant *B* is set to unity, the parameters responsible for the environment dynamics are set to  $\alpha=\beta=k$  – as before we will check the performance for different values of  $k=\{0.001, \dots, \infty\}$ 

(7)

1,1000}, the arrival rate of CRN traffic is  $\lambda$ =4. This set of parameters  $\Omega$  initializes the algorithm for CE based policy search described in Appendix B, the additional parameters controlling the algorithm are: population size *N*=1000, number of best samples *K*=10, maximum iterations *T*=100, threshold values *d*=5 and  $\varepsilon$ =1*e*-4.



Figure 10. The CE optimized policy for parameters  $M=6, \alpha=\beta=1e3, \mu=1, \lambda=4, C=(3,3,3,3,3,3,4)$  $\theta=(0.5653, 0.9125, 0.9254, 0.9900, 0.7184, 0.5883, 0.3729)$ 

In our associated stochastic problem the policy chooses action  $(C_t, \theta_t)$  when CRN is in state  $(\hat{S}_t, X_t)$ . We assume that,  $C_t$  is a discrete random variable that takes integer values  $\{0,1,\ldots,M\}$ , while the tradeoff parameter  $\theta_t$  is normally distributed according to a truncated normal distribution in the range [0,1]. Note that our policy is state dependent. We distinguish between the cases  $X_t=0$  and  $X_t>0$ . Obviously, for  $X_t=0$  CRN has no packets to transmit and in this case it is reasonable to allocate the bandwidth resources to the sensing process ( $\theta_t = 1$ ). The CE algorithm optimizes the policy for  $X_t>0$ .

The resulting CE optimized policies are presented in Figs. 10–12. For the case k=1000, presented in Fig.10, the environment changes are too fast for the perception process and CRN fails to keep track of the network state. This can be seen through the fact that the channel allocation C=(3,3,3,3,3,3,4) is insensitive to the estimation  $\hat{S}_t$ , and the number of accessed channels is approximately the average number of unoccupied channels. This is opposed to the cases with slower environment as we will see later. Nevertheless, the tradeoff parameter  $\theta=(0.5653,0.9125,0.9254,0.9900,0.7184,0.5883,0.3729)$  shows that CRN tries to avoid collisions with PU; it can be shown by a simple analysis of the CTMC (in Fig. 2) that for  $\alpha=\beta$ ,  $S_t$  resides only a small portion of time in the states 0 and M while it spends more time in the inner states. This fact is reflected in the low values of  $\theta$  when the perception process estimates the environment to be in states 0 or M. In this situation more resources are allocated to perception in order to better react to the fast transitions of the environment state.



Figure 11. The CE optimized policy for parameters  $M=6, \alpha=\beta=1, \mu=1, \lambda=4, C=(1,1,2,3,4,5,5)$  $\theta=(0.3284, 0.7125, 0.7769, 0.9152, 0.9243, 0.8718, 0.6919)$ 

In Fig. 11, the parameter k is set to 1. It can be seen that the resulting policy is more sensitive to the estimation of the environment state  $\hat{S}_t$ , and the number of accessed channels is approximately  $\hat{S}_t$  except for the rapidly switching states 0 and M. As in the previous case, the tradeoff parameter  $\theta$ , allocates more bandwidth for transmissions when the estimation  $\hat{S}_t$  indicates that the environment state is a persistent one, and it increases the sensing rate when the environment moves to quickly changing states.



Figure 12. The CE optimized policy for parameters  $M=6, \alpha=\beta=1e-3, \mu=1, \lambda=4, C=(1,1,2,3,4,5,6)$  $\theta=(0.9928, 0.9936, 0.9948, 0.9922, 0.9982, 0.9867, 0.9800)$ 

Finally, Fig. 12 represents situation where the environment changes occur in a significantly slower manner compared to the rate of the perception process. This fact can be observed from the tradeoff parameter  $\theta$ , which takes very high values independently of the estimation  $\hat{S}_t$ . The allocation of the channels *C* is equal to  $\hat{S}_t$  even for the rapidly switching state *M*.

Note that the optimized policies allow the sampling rate and the number of accessed channels to be a function of the current state estimation. This is crucial, because when some states are very likely to persist for longer periods, the cognition cycle may choose a more efficient course of action.

# V. SUMMARY

In this paper, a three-dimensional CTMC process has been introduced to model the operation of CRN where PU form a birth-death process and SU can queue. The analytical framework combines the environment dynamics, perception and decision making components of the cognition cycle and the spectrum access processes. The model was analyzed using matrix geometric approach. The analysis results give insights about the behavior of CRN in general and the impact of sensing rate and the system dynamics on the waiting times of secondary users in particular.

The cognition cycle is treated as an integral part of the system's overall behavior, and we optimize policies controlling simultaneously the interdependent perception and transmission processes. In this way, the resources are allocated according to the needs of the overall task. The CE optimized policies demonstrate adaptive behavior in which the resources are intelligently allocated to the perception and the transmission processes in a task-relevant manner.

# VI. APPENDIX A – ANALYSIS OF THE 3-D CTMC

In this appendix we present the analysis of the three dimensional CTMC illustrated in Fig. 4

# A. CTMC Structure

In order to make the analysis of the system easier we numerate the states of  $Z_t$  lexicographically, i.e.  $(0,0),(0,1),\ldots,(0,M),(1,0),(1,1),\ldots,(M,M)$  and index them 1 to  $(M+1)^2$ . This new order of states turns our CTMC to two dimensional since now  $Z_t \in \{1,2,\ldots,(M+1)^2\}$ . Then again we order the states lexicographically, i.e.  $(0,1),(0,2),\ldots,(0,M+1),(1,1),(1,2),\ldots$  and construct the generator matrix Q of this CTMC which is given by:

$$Q = \begin{pmatrix} B_{00} & B_{01} & 0 & 0 & 0 \\ B_{10} & B_{11} & A_0 & 0 & 0 \\ \hline 0 & A_2 & A_1 & A_0 & 0 & \cdots \\ 0 & 0 & A_2 & A_1 & A_0 \\ 0 & 0 & 0 & A_2 & A_1 \\ & \vdots & & \ddots \end{pmatrix}$$

where  $B_{00}=\{B_{00}(i,j)\}$ ,  $B_{01}=\{B_{01}(i,j)\}$ ,  $B_{10}=\{B_{10}(i,j)\}$ ,  $B_{11}=\{B_{11}(i,j)\}$ ,  $A_0=\{A_0(i,j)\}$ ,  $A_1=\{A_1(i,j)\}$  and  $A_2=\{A_2(i,j)\}$  are  $(M+1)^2\times(M+1)^2$  matrices. A 0 entry in Q (and in other matrices) is a matrix of all zeros of the appropriate dimension. It can be seen that in our model  $B_{01}=A_0=\text{diag}\{\lambda,\lambda,...,\lambda\}$ . For each value  $z_{i,j}=(i,j)$  the process  $Z_t$  can take, the service rate is  $\mu_{i,j}=\mu \min\{i,j\}$ . We order the elements  $\mu_{i,j}$  in the same way as we did for  $Z_t$  and obtain a vector of service rates  $\underline{\mu}$ . It can be seen that  $B_{10}=A_2=\text{diag}\{\underline{\mu}\}$ , while the matrices  $B_{00}$  and  $B_{11}=A_1$  are more complicated:

$$B_{00}(i,j) = \begin{cases} -\left(\lambda + (M - \lfloor i/(M+1) \rfloor)\alpha + \lceil i/(M+1) \rceil\beta\right) j = i \\ (M - \lfloor i/(M+1) \rfloor)\alpha & j = i + M \\ \lceil i/(M+1) \rceil\beta & j = i - M \\ \delta & j = \lfloor i/(M+1) \rfloor (M+2) \cap i \neq j \\ 0 & else \end{cases}$$

$$A_{1}(i,j) = \begin{cases} -\left(\lambda + i\mu + (M - \lfloor i/(M+1) \rfloor)\alpha + \lceil i/(M+1) \rceil\beta\right) j = i \\ (M - \lfloor i/(M+1) \rfloor)\alpha & j = i + M \\ \lceil i/(M+1) \rceil\beta & j = i - M \\ \delta & j = \lfloor i/(M+1) \rfloor (M+2) \cap i \neq j \\ 0 & else \end{cases}$$

#### **B.** Stationary Probabilities

We define the stationary probabilities  $\pi_{i,j}$  of the process to be at level *i* and state *j* within that level. Calculating the stationary probabilities will allow evaluating interesting quantities, mainly the waiting time of SU. The calculations here follow [24] and are adopted for our model.

Let  $\pi_i \equiv (\pi_{i,1}, \pi_{i,2}, \dots, \pi_{i,M^2})$  and  $\pi \equiv (\pi_0, \pi_1, \pi_2, \dots)$ . The stationary distribution is the unique set of  $\pi_i \ge 0$ , *i* $\ge 0$ , that solves

$$\begin{cases} \pi Q = \underline{0} \\ \pi \underline{e} = 1 \end{cases}$$
(A.1)

where  $\underline{e}$  (<u>0</u>) denotes an appropriately dimensioned column (row) vector of 1's (0's). From the first equation in (A.1) we may write down for the repeating portion of the process:

$$\pi_{j-1}A_0 + \pi_j A_1 + \pi_{j+1}A_2 = \underline{0} \qquad (j \ge 1)$$
(A.2)

For this type of CTMC characterized by a boundary conditions in the first column of Q followed by a repetitive portion of columns containing matrices  $A_0$ ,  $A_1$  and  $A_2$ , there exist some constant matrix R such that

$$\pi_j = \pi_{j-1} R, \qquad (j \ge 1) \tag{A.3}$$

and that the values of  $\pi_j$ ,  $j \ge 1$ , have a matrix geometric form, i.e.:

$$\pi_j = \pi_0 R^j, \quad (j \ge 1) \tag{A.4}$$

substituting (A.4) into (A.2) yields

$$A_0 + RA_1 + R^2 A_2 = 0 (A.5)$$

This quadratic equation in R is typically solved numerically. There is more than one R that solves (A.5). When the CTMC is ergodic, there is a unique stationary distribution  $\pi$  that satisfies (A.1). Analogous to the scalar case where the utilization factor should be less then unity, in our case all eigenvalues of R must be less then unity for the normalization constraint in (A.1) to hold [25].

After solving for R, in order to determine the stationary probabilities, we continue with the boundary conditions:

$$\pi_0 B_{00} + \pi_1 A_2 = \pi_0 (B_{00} + RA_2) = 0 \tag{A.6}$$

Equation (A.6) alone is not enough to solve for  $\pi_0$  since it is not of full rank and we must use the normalization constraint in (A.1):

$$\pi e = \left(\sum_{j=0}^{\infty} \pi_j\right) e = \pi_0 (I - R)^{-1} e = 1$$
(A.7)

Combining (A.6) and (A.7) we have

$$\pi_0[(I-R)^{-1}\underline{e}, (B_{00}+RA_2)^*] = [1,\underline{0}]$$
(A.8)

where  $(B_{00}+RA_2)^*$  is the result from removal of the first column from the matrix  $(B_{00}+RA_2)$ , and  $[1,\underline{0}]$  is a row vector consisting of a 1 followed by  $M^2-1$  zeros. Equation (A.8) is solved by appropriate numerical methods.

# VII. APPENDIX B – CROSS-ENTROPY ALGORITHM FOR CRN POLICY OPTIMIZATION

In this appendix, we present the CE algorithm for CRN policy optimization.

# Input:

- function  $W(P;\Omega)$
- system parameters  $\Omega = \{\alpha, \beta, M, \lambda, \mu\}$
- probability density families  $\{p_C(\cdot; \sigma_C)\}$  and  $\{p_\theta(\cdot; \sigma_\theta)\}$ ,
- initial parameters  $\sigma_{C,0}$  and  $\sigma_{\theta,0}$
- parameters N,K,T,d,  $\varepsilon$
- $\bullet \ t \leftarrow 0$

# Repeat

- 1: Generate samples  $C^{(k)}$  ( $k=1,2,\ldots,K$ ) from  $p_C(\cdot;\sigma_{C,t-1})$
- 2: Generate samples  $\theta^{(k)}$  (*k*=1,2,...,*K*) from  $p_{\theta}(\cdot;\sigma_{\theta,t-1})$
- 3: Compose policy samples  $P^{(k)} = (C^{(k)}, \theta^{(k)}) \ (k=1,2,...,K)$
- 4: Calculate  $W^{(k)} = W(P^{(k)}; \Omega)$  for each sample (k=1,2,...,K)
- 5: Keep  $N(N \le K)$  best samples graded by their  $W^{(k)}$  value and discard the other samples
- 6:  $V_t = \min_k(W^{(k)})$  (minimize over the saved *N* best samples)
- 7: Using the N best samples update the parameters

7.1: 
$$\sigma_{C_{d}} \leftarrow \arg\max_{\sigma_{C}} \sum_{n=1}^{N} \ln(p_{\sigma_{C}}(C^{(n)};\sigma_{C}))$$
  
7.2:  $\sigma_{\theta_{d}} \leftarrow \arg\max_{\sigma_{\theta}} \sum_{n=1}^{N} \ln(p_{\sigma_{\theta}}(\theta^{(n)};\sigma_{\theta}))$   
8:  $t \leftarrow t + 1$ 

**Until**  $(t > T \text{ or } |V_t - V_{t-\tau}| \le \varepsilon, \tau = 1, 2, ..., d)$ 

**<u>Output</u>**:  $P^* = (C^*, \theta^*)$  – best sample,  $W^* = W(P^*; \Omega)$  – best value

## VIII. REFERENCES

- [1] Youping Zhao, Shiwen Mao, James Neel, and Jeffrey Reed, "Performance Evaluation of Cognitive Radios: Metrics, Utility Functions and Methodologies," Proceedings of the IEEE vol 97, Issue 4, April 2009.
- [2] J. Mitola et al., "Cognitive radio: Making software radios more personal," IEEE Pers. Commun., vol. 6, no. 4, pp. 13–18, Aug. 1999.
- <sup>[3]</sup> J. Mitola, "Cognitive radio an integrated agent architecture for software defined radio," Ph.D. dissertation, Royal Institute of Technology, Kista, Sweden, May 8 2000.
- [4] S. Haykin, "Cognitive radio: brain-empowered wireless communications," IEEE J. on Selected Areas in Communications, vol. 23 pp. 201-220, February, 2005.
- [5] Boyd, J. "A discourse on winning and losing," Maxwell Air Force Base, AL: Air University Library Document No. M-U 43947 (Briefing slides) ,1987.
- [6] I. F. Akyildiz, L. Won-Yeol, M. C. Vuran, and S. Mohanty, "A survey on spectrum management in cognitive radio networks," IEEE Communications Magazine, vol. 46, no. 4, pp. 40–48, 2008.
- [7] Federal Communications Commission, ET Docket No 03-222 Notice of proposed rule making and order, December 2003.
- [8] Federal Communications Commission. Spectrum Policy Task Force Report. ET Docket No. 02-135, November 2002.
- [9] R.W. Brodersen, A. Wolisz, D. Cabric, S.M. Mishra, D. Willkomm, Corvus: "A cognitive radio approach for usage of virtual unlicensed spectrum," Berkeley Wireless Research Center (BWRC) White paper, 2004.
- [10] FCC, "Facilitating Opportunities for Flexible, Efficient, and Reliable Spectrum Use Employing Cognitive Radio Technologies," Notice of Proposed Rule Making and Order, Docket No. 03-108, Dec 30, 2003.
- [11] I. J. Mitola, "Software radios: Survey, critical evaluation and future directions," IEEE Aerosp.Electron. Syst. Mag, vol.8, pp.25-36, Apr. 1993.
- [12] I. F. Akyildiz, W.-Y. Lee, M. C. Vuran, and S. Mohanty, "NeXt generation/dynamic spectrum access/cognitive radio wireless networks: a survey," Computer Networks, vol.50 pp.2127-2159, 2006.
- [13] A. T. Hoang and Y. C. Liang, "Adaptive Scheduling of Spectrum Sensing Periods in Cognitive Radio Networks", in Proc. IEEE GlobeCom 2007, Washington D.C., USA
- [14] Y.-C. Liang, Y. Zeng, E. Peh, and A. T. Hoang, "Sensing-throughput tradeoff for cognitive radio networks," in Proc. IEEE Int. Conf. Commun.(ICC), June 2006, pp. 5330-5335.
- [15] Y. Pei, A. T. Hoang, and Y.-C. Liang, "Sensing-throughput tradeoff in cognitive radio networks: how frequently should spectrum sensing be carried out?" in Proceedings of the 18th Annual IEEE International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC '07), Athens, Greece, September 2007.

- [16] A. Ghasemi and E. S. Sousa, "Optimization of spectrum sensing for opportunistic spectrum access in cognitive radio networks," in Proc. 4th IEEE Consumer Commun. Networking Conf. (CCNC), Jan. 2007, pp.1022-1026.
- [17] C. Cormio and K. R. Chowdhury, "A survey on MAC protocols for cognitive radio networks," Ad Hoc Networks, vol. 7, no. 7, pp. 1315-1329, September 2009.
- <sup>[18]</sup> D. Maldonado, B. Le, A. Hugine, T.W. Rondeau, C.W. Bostian, "Cognitive radio applications to dynamic spectrum allocation: a discussion and an illustrative example", DySPAN, Nov. 2005.
- [19] Shu, T. and Krunz, M. 2009. Throughput-efficient sequential channel sensing and probing in cognitive radio networks under sensing errors. In *Proceedings of the 15th Annual international Conference on Mobile Computing and Networking*(Beijing, China, September 20 - 25, 2009). MobiCom '09. ACM, New York, NY, 37-48.
- [20] Rashid, M. M., Hossain, M. J., Hossain, E., and Bhargava, V. K., "Opportunistic spectrum scheduling for multiuser cognitive radio: a queueing analysis". *Trans. Wireless. Comm.* 8, 10 (Oct. 2009), 5259-5269.
- [21] B.Oklander, M.Sidi, "Modeling and Analysis of System Dynamics and State Estimation in Cognitive Radio Networks", PIMRC'10 (CogCloud Workshop), September 2010, Istanbul, Turkey.
- [22] M.F. Neuts, "Matrix Geometric Solutions in Stochastic Models", John Hopkins University Press, 1981.
- [23] T.A. Weiss and F.K. Jondral, "Spectrum pooling: an innovative strategy for the enhancement of spectrum efficiency," IEEE Commun. Mag., vol. 42, no. 3, pp. 8–14, Mar. 2004.
- <sup>[24]</sup> R. Nelson, "Matrix Geometric Solutions in Markov Models: A Mathematical Tutorial," IBM Research Report RC 16777, 1991.
- [25] R.A. Horn, C.R. Johnson, "Matrix Analysis", Cambridge, 1987.
- [26] R. Y. Rubinstein and D. P. Kroese, The Cross Entropy Method. A Unified Approach to Combinatorial Optimization, Monte-Carlo Simulation, and Machine Learning, ser. Information Science and Statistics, M. Jordan, J. Kleinberg, B. Scholkopf, F. Kelly, and I. Witten, Eds. Springer, 2004.
- [27] S. Mannor, R. Y. Rubinstein, and Y. Gat, "The cross-entropy method for fast policy search," in Proceedings 20th International Conference on Machine Learning (ICML-03), Washington, US, 21–24 August 2003,
- [28] L. Busoniu, D. Ernst, B. De Schutter, and R. Babuska, "Policy search with crossentropy optimization of basis functions," *Proceedings of the 2009 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL 2009)*, Nashville, Tennessee, pp. 153-160, Mar.-Apr. 2009.
- [29] J. Mitola III, Cognitive Radio: An Integrated Agent Architecture for Software Defined Radio, PhD Thesis, Royal Institute of Technology (KTH), Sweden, 8 May, 2000.