Competing Cognitive Resilient Networks

Siamak Dastangoo, Member, IEEE, Carl E. Fossa, Senior Member, IEEE, Youngjune L. Gwon, Member, IEEE, and Hsiang-Tsung Kung

Abstract—We introduce competing cognitive resilient 1 2 network (CCRN) of mobile radios challenged to optimize data 3 throughput and networking efficiency under dynamic spectrum 4 access and adversarial threats (e.g., jamming). Unlike the 5 conventional approaches, CCRN features both communicator 6 and jamming nodes in a friendly coalition to take joint actions 7 against hostile networking entities. In particular, this paper 8 showcases hypothetical blue force and red force CCRNs and 9 their competition for open spectrum resources. We present state-10 agnostic and stateful solution approaches based on the decision 11 theoretic framework. The state-agnostic approach builds on 12 multiarmed bandit to develop an optimal strategy that enables 13 the exploratory-exploitative actions from sequential sampling 14 of channel rewards. The stateful approach makes an explicit 15 model of states and actions from an underlying Markov decision 16 process and uses multiagent *Q*-learning to compute optimal 17 node actions. We provide a theoretical framework for CCRN 18 and propose new algorithms for both approaches. Simulation 19 results indicate that the proposed algorithms outperform some 20 of the most important algorithms known to date.

21 *Index Terms*—Cognitive radio, strategy, multi-armed ban-22 dit (MAB), reinforcement learning.

23

I. INTRODUCTION

²⁴ C OGNITIVE radios have arisen commercially over the last ²⁵ C decade, enabling a new means to share radio spectrum. ²⁶ Dynamic spectrum access (DSA) [1] is a compelling usage ²⁷ scenario for cognitive radio systems. DSA aims to relieve ²⁸ shortages of radio spectrum, which is the scarcest—hence, the ²⁹ most expensive—resource to build a wireless network. Much ³⁰ of contemporary research has considered cognitive radios as ³¹ the secondary user of a licensed spectrum and focused on ³² the development of a flexible mechanism to opportunistically ³³ access the licensed channel to its maximal spectral efficiency. ³⁴ We envision the use of cognitive radio technology for tac-

³⁴ We envision the use of cognitive radio technology for the ³⁵ tical wireless networks operating in an environment where

Manuscript received October 23, 2015; revised March 8, 2016; accepted April 18, 2016. This material is based upon work supported by the Office of the Secretary of Defense under Air Force Contract FA8721-05-C-0002 and/or FA8702-15-D-0001. The associate editor coordinating the review of this paper and approving it for publication was C. Clancy.

S. Dastangoo and C. E. Fossa are with the Tactical Networks Group, MIT Lincoln Laboratory, Lexington, MA 02421 USA (e-mail: sia@ll.mit.edu; cfossa@ll.mit.edu).

Y. L. Gwon was with Computer Science at Harvard John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138 USA. He is now with MIT Lincoln Laboratory, Lexington, MA 02421 USA (e-mail: gyj@ll.mit.edu).

H.-T. Kung is the William H. Gates Professor of Computer Science and Electrical Engineering at Harvard John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138 USA (e-mail: kung@harvard.edu).

Digital Object Identifier 10.1109/TCCN.2016.2570798

malicious jammers and other adversarial threats exist in addition to scarce spectral resources. The past approaches [2]–[6] ³⁷ have concentrated on a defense mechanism against the adversarial jamming attacks (i.e., antijamming strategy). In contrast, our approach in this paper is to optimize antijamming and jamming actions jointly. ⁴¹

This paper introduces Competing Cognitive Resilient 42 Network (CCRN), where a network of communicator (comm) 43 nodes and jammers attempts to dominate the access to an open 44 spectrum *competitively* against a hostile opponent, which is 45 possibly another cognitive radio network of similar capabil-46 ities. An antijamming-jamming strategy is a critical require-47 ment for CCRNs, and we propose two different approaches in 48 computing the optimal CCRN strategy, namely state-agnostic 49 and stateful. 50

The state-agnostic approach is based on Multi-armed 51 Bandit (MAB) problems [7] that address the exploration-52 exploitation dilemma for allocating resources on sequential 53 reward sampling. Lai et al. [8], for example, discusses the 54 primary-secondary usage framework in the context of MAB 55 for cognitive radios in the DSA paradigm. In our work, 56 however, we have devised a randomized algorithm for the 57 CCRN nodes taking actions (i.e., communicate or jam) on a 58 block of multi-channel spectrum that works statelessly and is 59 guided only by channel reward sampling. The state-agnostic 60 algorithm essentially runs Thompson sampling [9], an old 61 probability matching heuristic, under which we set up the optimal Bayesian conjugate prior from an extreme-valued like-63 lihood. We will explain the detailed rationale and present 64 its superior performance over some of the most important 65 MAB algorithms applied to CCRNs in later sections of this 66 paper. 67

The stateful approach, on the other hand, is based on game 68 theory [10]. It requires explicit modeling of CCRN states and 69 actions from an underlying Markov Decision Process (MDP). 70 We have formulated the two-network CCRN game and decom-71 pose it to antijamming and jamming subgames. Unlike the 72 existing game-theoretic framework for cognitive radios, we solve the antijamming and jamming subgames jointly for an 74 optimal strategy, applying multi-agent Q-learning [11]. Given 75 perfect sensing at the lower layer, we will show that Q-learning 76 can result in channel access decisions that converge to the best 77 cumulative average reward in the steady state. 78

We summarize our main contributions:

This paper develops a theoretical framework for cognitive radio networks under competition that simultaneously access (communicate) and suppress (jam) for spectral dominance in tactical setting;

79

2332-7731 © 2016 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

- We propose new algorithmic approaches to compute a
- ⁸⁵ joint antijamming-jamming strategy in-network;
- We evaluate the empirical performance of the proposed
- algorithms and show that they outperform the key classi-cal approaches known to date.

This paper is an extended version of our two previous pub-89 ⁹⁰ lications in tactical cognitive radio networks [12], [13]. The ⁹¹ paper provides complete system and architectural consider-92 ations for an intelligent tactical network. In particular, we 93 describe a detailed architecture under which cognitive radio 94 nodes collect information by sensing and assess the perfor-95 mance of learning-based algorithms for dyanamic spectrum ⁹⁶ access. Furthermore, we validate the optimality of our algo-⁹⁷ rithmic approaches and their equivalence to analytical models, which are generally known intractable, using probabilistic 98 99 sampling (i.e., by Markov Chaning Monte Carlo method) on 100 the derived transition probability functions. This paper also ¹⁰¹ presents the complete reward performance results by iterating ¹⁰² all comm and jamming probabilities of cognitive radio nodes. The rest of the paper is organized as follows. In Section II, 103 104 we describe the CCRN system architecture. Section III 105 presents our mathematical formulation. Section IV develops ¹⁰⁶ algorithmic approaches to find an optimal CCRN strategy. We describe our adaptation of classical algorithms under the 107 108 state-agnostic and stateful approaches and propose two new ¹⁰⁹ algorithms. In Section V, we numerically evaluate the perfor-110 mance of the proposed algorithms. Section VI discusses related ¹¹¹ work, and Section VI concludes the paper.

II. COMPETING COGNITIVE RESILIENT NETWORK (CCRN): SYSTEM ARCHITECTURE

114 A. Overview

For clarity of discussion, imagine two networks of cognitive radios, Blue Force (BF or the *ally*) and Red Force (RF or the *enemy*). Each network consists of two types of nodes: communicator (comm) and jammer. In their field operations, the BF and RF networks face each other in a competition to achieve higher comm data throughput and suppress the oppotern nent's comm activities by jamming, all trying to dominate an tree open spectrum.

Accessing channels by a comm node is determined upon sensing and cognizance of vacant spectrum blocks. The primary-secondary user dichotomy popular in the existing cognitive radio networking literature is mostly invalid for our purpose. We use the term Competing Cognitive Resilient Network Resilient Network to designate our BF and RF networks instead.

In tactical ad hoc networks, node geometry in the field can be critical for effective antijamming and jamming capabilities. For simplicity, we assume that the internal layout of the nodes for BF or RF networks does not play a significant role. However, we use two different network control mechanisms, namely centralized and distributed. The centralized control for CCRN designates a central entity that collects channel sensing (performed by each node) and other network information to make a coherent, network-wide decision to access or suppress a channel. For each time slot, the centralized control not provide the channels to broadcast control information.



Fig. 1. Transmission opportunity $\langle f_i, B_i, t, T \rangle$ (shaded region).

The information on control channel selection is transmitted 140 through the previous slot's control channel. Under the distributed control, each node in a CCRN acts as a decision maker 142 for its own action after exchanging the information with other 143 nodes. Due to staleness of information at each node, however, 144 distributed decision making may result in suboptimal actions 145 for the network as a whole. 146

B. Communications Model

Spectrum for open access is partitioned in time and frequency. There are *N* non-overlapping channels located at ¹⁴⁹ the center frequency f_i (MHz) with bandwidth B_i (Hz) for ¹⁵⁰ i = 1, ..., N. A transmission opportunity is represented by a ¹⁵¹ tuple $\langle f_i, B_i, t, T \rangle$, which designates a time-frequency slot at ¹⁵² channel *i* and time *t* with time duration *T* (msec) as depicted ¹⁵³ in Fig. 1. We assume a simple CSMA in which comm nodes ¹⁵⁴ first sense before transmitting in a slot of opportunity. ¹⁵⁵

In order to coordinate a non-conflicting spectrum access and 156 jamming strategy network-wide, we assume that the nodes 157 (both comm and jammers) exchange necessary information 158 via control messages. We call the channels used to exchange 159 control messages 'control channels.' On the contrary, 'data 160 channels' are used to transport regular data packets. We follow 161 the DSA approach [2] that control or data channels are dynamically allocated. Due to cryptographic randomization, guessing 163 the correct control channel by the enemy network would be a 164 hard problem. However, if the control channel happens to be 165 jammed, say at time *t*, the spectrum access at time t + 1 will 166 be uncoordinated. For the case of blocked control channels, 167 the CCRN performance will be suboptimal.

C. Jamming Model

Xu *et al.* [14] presents a sound taxonomy of RF jamming. A constant jammer continuously dissipates power into a selected channel by transmitting arbitrary waveforms. A deceptive jammer can instead send junk bits encapsulated in a legitimate packet to conceal its intent to disrupt comm nodes. A random jammer alternates between jamming and remaining quiet for random time intervals. A reactive jammer listens to a channel, stays quiet when the channel is idle, and starts transmitting upon sensing an activity.

The key to successful jamming is intelligence and ¹⁷⁹ statistical sophistication. We introduce *strategic* jamming ¹⁸⁰ model that extends the statistical jamming described in ¹⁸¹ Pajic and Mangharam [15]. A strategic jammer can learn ¹⁸² media access patterns of comm nodes and adapt to antijam-¹⁸³ ming schemes in the same way that our comm nodes leverage ¹⁸⁴ sensing and cognition. Strategic jammers can remain effective ¹⁸⁵ for long without being detected, subsequently causing more ¹⁸⁶ damages than the jammers based on existing models. ¹⁸⁷

169

147

TABLE I NODE ACTIONS, OUTCOME AND RESULTING REWARD

BF	BF	RF	RF		
comm	jammer	comm	jammer	Outcome	Reward
Тx	Ø	Ø	Ø	BF Tx success	$R_B += B$
Ø	Jam	Tx	Ø	BF jamming	$R_B += B$
Tx	Jam	Ø	Ø	BF misjamming	-
Ø	Ø	Τx	Ø	RF Tx success	$R_R += B$
Tx	Ø	Ø	Jam	RF jamming	$R_R += B$
Ø	Ø	Tx	Jam	RF misjamming	-
Tx	Ø	Tx	Ø	Tx collision	-

188 D. Reward Model

We employ a reward metric measured in *bits*. When a CCRN response to the makes successful transmission of a packet containing *B* bits of data, it receives the reward of *B* (bits). The definition of a successful transmission follows the notion in response to the reward of *B* (bits) and the matrix classical wireless networking that there should be only one response to the transmission for the Tx opportunity. If there were response two or more simultaneous comm transmissions (from either response to the same or different network), a collision occurs, and no response to the transmission for the transmission occurs, and no response to the same or different network).

Jammers do not create any reward by themselves. However, ¹⁹⁹ they receive a reward by suppressing an opposing comm ²⁰⁰ node's otherwise successful transmission. For example, a BF ²⁰¹ jammer earns a reward *B* by jamming the slot in which a sole ²⁰² RF comm node transmits. If there were no jamming, the RF ²⁰³ comm node would have earned *B*. Also, a BF jammer can ²⁰⁴ jam a BF comm mistakenly (e.g., due to faulty intra-network ²⁰⁵ coordination), which we call *misjamming*.

205 coordination), which we call *misjamming*. Table I summarizes how the outcome and reward at a chan-206 207 nel are determined given various combinations of the BF and 208 RF node actions at a slot of transmission opportunity ('Ø' 209 means no action). Note that each CCRN keeps track of its ²¹⁰ reward cumulatively (i.e., R_B , R_R) over operated time slots. For illustrative purposes, we provide example BF and RF 211 ²¹² actions and explain the reward computation for each CCRN. 213 Let BF and RF networks each have two comm nodes and two ²¹⁴ jammers. At time t, BF comm node 1 transmits in channel 7, 215 and BF comm node 2 in channel 3. BF jammers jam channels 1 and 5 at t. Similarly, RF has its comm nodes transmitting in 216 channels 3 and 5, and its jammers jamming channels 10 and 9. 217 218 Fig. 2 depicts the resulting channel-action bitmap where 1 219 indicates transmit or jam and 0 otherwise. BF jammer on chan-220 nel 5 is successful whereas the one on channel 1 is not. There is a comm collision between BF and RF in channel 3, but BF 221 222 has a successful comm transmission in channel 7. Thus, BF 223 network receives the reward of 2B for one of its comm and 224 one of its jamming actions. RF networks has no success in 225 comm or jamming at t.

226 E. System Model

We now describe a CCRN system in detail. A CCRN node consists of sensing, strategy, schedule, and Tx/jam components as illustrated in Fig. 3. Using local and global sensing information, a CCRN node applies a strategy to compute an action (i.e., transmit, jam, or do nothing) particular to its channel of



Fig. 2. Example Blue Force (BF) and Red Force (RF) CCRN node actions and resulting outcome.



Fig. 3. Overview of Competing Cognitive Resilient Network (CCRN).

interest. The action translates to scheduling a block of transmission/jamming opportunities. Node actions can be computed in either centralized or distributed manner.

Under the centralized control, CCRN works in the following	235
steps.	236
1) Sense channel activities (each node)	237
2) Collect sensing information (controller)	238
3) Compute node actions (controller)	239
4) Disseminate node actions (controller)	240
5) Act on channel (each node)	241
In the distributed control, CCRN works as follows.	242
1) Sense channel activities (each node)	243
2) Exchange sensing information (each node)	244
3) Compute its own action (each node)	245
4) Act on channel (each node)	246
When acting on a channel, a CCRN node under the dis-	247

when acting on a channel, a CCRN node under the dis-247 tributed control assumes *cooperative* behavior by considering 248 the holistic goal of maximizing the overall network performance. When new information is available for a node (e.g., 250 channel sensing, outcome), a CCRN node exchanges the 251 information *collaboratively* with one another. The distributed 252 CCRN would be a realistic application for Mobile ad hoc 253 network (MANET) where there is little or no fixed infrastructural support. In the later sections of this paper, we will 255 evaluate CCRN under both centralized and distributed control 256 architectures. 257

III. MATHEMATICAL FORMULATION 258

This section presents mathematical formulation for CCRN. ²⁵⁹ In particular, we develop two compatible, yet different ²⁶⁰

TA	BL	ΕII
UMMARY	OF	NOTATIONS

S

Notation	Description
a_B^t, a_R^t	node actions at t for Blue (B) and Red (R) Force networks
Γ^t	regret over t time slots
R^t	cumulative reward over t time slots
r^{j}	instantaneous reward at time slot j
γ	reward discount ratio
σ^t	channel access strategy at time t
π	probability distribution computed over action set
Ω^t	outcome bit vector (each 1 is success, and 0 failure) at t
heta, a, b, eta	hyperparameters of probability model
Q(s,a)	measure for quality of action at state s
V(s)	measure for value of state
$lpha,\lambda$	learning rate, decay parameters for Q-learning
N	total number of accessible channels
C, J	number of comm and jamming nodes
p_{Tx}	transmit probability of a comm node

²⁶¹ frameworks for our BF-RF competition scenario. First, we ²⁶² explain a *state-agnostic* model that operates a CCRN without ²⁶³ considering any system states. Then, we describe a *stateful* ²⁶⁴ model by explicitly defining CCRN states and discuss plausi-²⁶⁵ ble ways in computing them. Markov Decision Process (MDP) ²⁶⁶ underlies the stateful CCRN whereas the state-agnostic coun-²⁶⁷ terpart is solely driven by sequential reward sampling.

268 A. Notation and Preliminaries

Strategy or policy for a CCRN means rules to select its 269 $_{270}$ node actions. There are C comm nodes and J jammers in $_{271}$ each CCRN. Let N designate the number of channels in the 272 spectrum. We define the BF action set A_B such that a BF ²⁷³ action $a_B \in A_B$, and similarly for RF, $a_R \in A_R$. At time a_{274} t, the BF and RF actions are $a_B^t = \{a_{B,comm}^t, a_{B,jam}^t\}$ and $a_{R}^{t} = \{a_{R,comm}^{t}, a_{R,jam}^{t}\}$ containing both comm and jamming actions, the size-*C* vectors $a_{B,comm}^t$, $a_{R,comm}^t$ and the size-*J* $a_{B,jam}^t$, $a_{R,jam}^t$. (Note we use a superscripted t for 'at time *t*.') 278 An *i*th element in $a_{B,comm}^t$ designates the channel number that 279 the *i*th BF comm node tries to transmit at *t*. Similarly, a *j*th 280 element in $a_{B,iam}^{t}$ is the channel that the *j*th BF jammer tries ²⁸¹ to jam at *t*. The BF and RF node actions result in an outcome $_{282} \Omega : A_B \times A_R \longrightarrow \mathbb{R}^N$. Subsequently, we can map the out-283 come to a reward $R: \Omega \longrightarrow \mathbb{R}$ given a reward function. In ²⁸⁴ Table II, we summarize a list of important notations used in 285 this paper.

286 B. State-Agnostic CCRN

29

The state-agnostic CCRN model is based on multi-armed bandit (MAB) originated from Thompson's medical experiments [9], although MAB is best explained with a gambler facing N slot machines (arms). The gambler's objective is to find a strategy that maximizes $R^t = \sum_{j=1}^{t} r^j$ for some t, the *cumulative* reward over a finite horizon. Lai and Robbins [7] introduced the concept of *regret* for a strategy σ measuring the distance from optimality

$$\Gamma^t = t\mu^* - \mathbb{E}\big[R^t_\sigma\big]$$

²⁹⁶ where μ^* is the hypothetical, maximum average reward if ²⁹⁷ gambler's action led to the best possible outcome each round ²⁹⁸ and achieved the actual reward R^t_{σ} under σ . It turns out that minimizing Γ^t is mathematically more convenient than ²⁹⁹ maximizing the *expectation* of R^t .

An arm corresponds to a channel in the spectrum under ³⁰¹ competition. Comm nodes and jammers are the players that a ³⁰² CCRN allocates to play (i.e., transmit or jam) the channels. ³⁰³ Since the CCRN has multiple nodes, our problem is classified ³⁰⁴ as multi-player MAB [16], which is different from the classic ³⁰⁵ single-player MAB formulated by Lai and Robbins [7]. In ³⁰⁶ addition, we have two system variations depending on whether ³⁰⁷ a centralized or distributed control mechanism is deployed. ³⁰⁸

Lai and Robbins [7] further derived the mathematical 309 qualification for an optimal strategy: 310

$$\lim_{t \to \infty} \sup \mathbb{E} \left[T_i^t \right] \le \frac{\log t}{\mathsf{D}_{KL}(p_i \parallel p^*)} \tag{2} \quad 311$$

where sup means supremum, T_i^I is the total number for arm ³¹² *i* being played, and $D_{KL}(\cdot \| \cdot)$ is the Kullback-Leibler diver-³¹³ gence [17] measuring the dissimilarity between the probability ³¹⁴ distributions p_i and p^* for the *i*-th arm's reward and the max-³¹⁵ imum reward resulted by choosing only the best possible arm ³¹⁶ each time. Eq. (2) provides the least upper bound for the num-³¹⁷ ber of times should an optimal arm—which could be different ³¹⁸ each time—be played asymptotically. Lai and Robbins also ³¹⁹ provided an algorithm that satisfies the condition of Eq. (2), ³²⁰ which will be discussed in Section IV. ³²¹

The BF strategy σ_B^I is a function over time. It takes ³²² necessary information such as sensing results and past actionoutcome/reward statistics as input and determines the BF node ³²⁴ actions. Under the centralized decision making, we express ³²⁵

$$\left\{x_B^j\right\}_{j=1}^t, \left\{a_B^j, \Omega^j\right\}_{j=1}^{t-1} \xrightarrow{\sigma_B^i} a_B^t \tag{3} 326$$

327

where x_B^t is the BF sensing results at t.

х

(1)

Under the distributed decision making, each node in the $_{328}$ network computes its own action. For BF node *i* (whether it $_{329}$ is a comm node or jammer), we write $_{330}$

$${}^{t}_{B,i}, \left\{ x_{B}^{j}, a_{B}^{j}, \Omega^{j} \right\}_{j=1}^{t-1} \xrightarrow{\sigma_{B,i}^{t}} a_{B,i}^{t}$$

$$(4) \quad {}^{331}$$

where $x_{B,i}^{t}$ is the sensing information only available to BF ³³² node *i* at time *t*, and $\sigma_{B,i}^{t}$ the strategy of BF node *i*'s own. At ³³³ time *t*, BF node *i* does not yet have all sensing results except ³³⁴ its own $x_{B,i}^{t}$. For the distributed case, node strategies can differ, ³³⁵ and there is no guarantee that conflicting actions of the nodes ³³⁶ in the same network such as collision and misjamming are ³³⁷ resolved. ³³⁸

C. Stateful CCRN 339

The stateful CCRN model uses the tuple $\langle S, A_B, A_R, R, T \rangle$ ³⁴⁰ to describe the dynamics of the competition between BF ³⁴¹ and RF networks. *S* denotes the state set, and $A_B =$ ³⁴² $\{A_{B,comm}, A_{B,jam}\}, A_R = \{A_{R,comm}, A_{R,jam}\}$ are the action sets ³⁴³ for BF and RF networks. The reward function $R : S \times$ ³⁴⁴ $\prod A_{\{B,R\},\{comm,jam\}} \rightarrow \mathbb{R}$ maps CCRN node actions to a ³⁴⁵ reward value at a given state. The state transition $T : S \times$ ³⁴⁶ $\prod A_{\{B,R\},\{comm,jam\}} \rightarrow PD(S)$ is the probability distribution ³⁴⁷ over *S*. Under stochastic setup, a strategy $\pi : S \rightarrow PD(A)$ ³⁴⁸ is the probability distribution over the action set. ³⁴⁹

We formulate the stateful CCRN with a Markov game [10]. 350

TABLE III COLLISION PARAMETERS

Parameter	Description
$I_{B,C}$	# of control channel collisions caused by BF comms only
$I_{B,D}$	# of data channel collisions caused by BF comms only
$I_{R,C}$	# of control channel collisions caused by RF comms only
$I_{R,D}$	# of data channel collisions caused by RF comms only
$I_{BR,C}$	# of control channel collisions caused by BF and RF comms
$I_{BR,D}$	# of data channel collisions caused by BF and RF comms

TABLE IV JAMMING PARAMETERS

Parameter	Description
$J_{B,R,C}$	# of BF control channel jammed by RF jammers
$J_{B,R,D}$	# of BF data channel jammed by RF jammers
$J_{B,B,C}$	# of BF control channel jammed by BF jammers
$J_{B,B,D}$	# of BF data channel jammed by BF jammers
$J_{B,BR,C}$	# of BF control channel jammed by BF and RF jammers
$J_{B,BR,D}$	# of BF data channel jammed by BF and RF jammers
$J_{R,B,C}$	# of RF control channel jammed by BF jammers
$J_{R,B,D}$	# of RF data channel jammed by BF jammers
$J_{R,R,C}$	# of RF control channel jammed by RF jammers
$J_{R,R,D}$	# of RF data channel jammed by RF jammers
$J_{R,BR,C}$	# of RF control channel jammed by BF and RF jammers
$J_{R,BR,D}$	# of RF data channel jammed by BF and RF jammers

1) State Representation: Each of N channels in the spec-351 352 trum under competition is described by a Markov chain. If there are L discrete states for a channel, we require to track L^N 353 states to describe the CCRN interactions fully. Such tracking, 354 however, leads to a computational complexity class in $O(L^N)$ 355 with respect to channelization N. We instead choose a terser 356 state representation $s = \langle I_C, I_D, J_C, J_D \rangle$ where I_C denotes the 358 number of control channels collided, I_D the number of data $_{359}$ channels collided, J_C the number of control channels jammed, and J_D the number of data channels jammed.

Given the current state and the action sets of BF and RF 361 ³⁶² nodes, the next state of CCRN is computable. The actions ³⁶³ of the opponent is inferred from channel measurements and sensing. To estimate I_C , I_D , J_C , and J_D , we need to observe 365 the parameters in Tables III and IV to calculate

 $I_C = \sum_{x \in \{B,R,BR\}} I_{x,C}$

 $I_D = \sum_{x \in \{B, R, BR\}} I_{x, D}$ $I_C = \sum_{x \in \{B, R, BR\}} I_{x, D}$

366

369

$$J_D = \sum_{x \in \{B,R\}, y \in \{B,R,BR\}} J_{x,y,D}$$

 $J_C =$

2) State Transition Probabilities: In this section, we derive 370 the full, analytical formula for the CCRN state transition 371 probability distribution that can be used for numerical approx-372 373 imation.

 $J_{x.v.c}$

a) Counting parameters for state transition: The fol-374 375 lowing conditional probability distribution determines the 376 transition function T:

5

To express I_C^{t+1} , I_D^{t+1} , J_C^{t+1} , and J_D^{t+1} , we need to define the ³⁷⁹ counting parameters related to collision and jamming: ³⁸⁰

- $m_{C1} \stackrel{\text{def}}{=} \#$ of collided control channels previously uncol- 381 lided and unjammed; 382
- $m_{C2} \stackrel{\text{def}}{=} \#$ of collided control channels previously collided; 383
- $m_{C3} \stackrel{\text{def}}{=} \#$ of collided control channels previously jammed; 384
- $m_{D1} \stackrel{\text{def}}{=} \#$ of collided data channels previously uncollided 385 and unjammed;
- $m_{D2} \stackrel{\text{def}}{=} \#$ of collided data channels previously collided; 387
- $m_{D3} \stackrel{\text{def}}{=} \#$ of collided data channels previously jammed; 388
- $n_{C1} \stackrel{\text{def}}{=} \#$ of jammed control channels previously uncol- 389 lided and unjammed; 390
- $n_{C2} \stackrel{\text{def}}{=} \#$ of jammed control channels previously collided; 391
- $n_{C3} \stackrel{\text{def}}{=} \#$ of jammed control channels previously jammed; 392
- $n_{D1} \stackrel{\text{def}}{=} \#$ of jammed data channels previously uncollided 393 and unjammed; 394
- 395

n_{D2} def # of jammed data channels previously collided;
n_{D3} def # of jammed data channels previously jammed. 206

Now we can write the number of collided control channels 397 $I_C^{t+1} = m_{C1} + m_{C2} + m_{C3}$, the total number of collided data 398 channels $I_D^{t+1} = m_{D1} + m_{D2} + m_{D3}$, the jammed control chan- 399 nels $J_C^{t+1} = n_{C1} + n_{C2} + n_{C3}$, and the jammed data channels 400 $J_D^{t+1} = n_{D1} + n_{D2} + n_{D3}.$ 401

We define the counting parameters that describe how BF 402 and RF networks choose control and data channels at time t: 403

- $\alpha_{C1}^{t} \stackrel{\text{def}}{=} \#$ of control channels chosen from previously 404 uncollided and unjammed channel space; 405
- $\alpha_{D1}^{t} \stackrel{\text{def}}{=} \#$ of data channels chosen from previously 406 uncollided and unjammed channel space; 407
- $\alpha_{C2}^{t} \stackrel{\text{def}}{=} \#$ of control channels chosen from previously 408 collided channel space; 409
- $\alpha_{D2}^t \stackrel{\text{def}}{=} \#$ of data channels chosen from previously 410 collided channel space;
- $\alpha_{C3}^t \stackrel{\text{def}}{=} \#$ of control channels chosen from previously $_{412}$ jammed channel space; 413
- $\alpha_{D3}^t \stackrel{\text{def}}{=} \#$ of data channels chosen from previously 414 jammed channel space. 415

We define the parameters to describe how BF and RF 416 jamming actions are chosen at t: 417

- $\alpha_{11}^{t} \stackrel{\text{def}}{=} \#$ of channels chosen from previously uncollided 418 channel space for jamming;
- $\alpha_{12}^t \stackrel{\text{def}}{=} \#$ of channels chosen from previously collided 420 channel space for jamming; 421
- $\alpha_{J1}^{t} \stackrel{\text{def}}{=} \#$ of channels chosen from previously unjammed 422 channel space for jamming; 423
- $\alpha_{I2}^t \stackrel{\text{def}}{=} \#$ of channels chosen from previously jammed 424 channel space for jamming. 425

We have a constraint $\alpha_{C1}^t + \alpha_{D1}^t < N_1^t$ where $N_1^t = N - {}_{426}^t$ $(I_C^t + I_D^t + J_C^t + J_D^t)$ gives the total number of uncollided and ${}_{427}^t$ unjammed channels. We also have $\alpha_{C2}^t + \alpha_{D2}^t < N_2^t$ where $_{428}^{428}$ $N_2^t = I_C^t + I_D^t$ is the total number of collided channels, and $_{429}^{429}$ $\alpha_{C3}^t + \alpha_{D3}^t < N_3^t$ where $N_3^t = J_C^t + J_D^t$ is the total number of $_{430}^{430}$ jammed channels. 431

b) Combinatorial analysis: We should consider combitant $(m_{C\{1,2,3\}}, m_{D\{1,2,3\}})$ and $(n_{C\{1,2,3\}}, n_{D\{1,2,3\}})$ subject to the constraints represented by I_C , I_D , J_C , and J_D . Using the binomial coefficient $\binom{n}{k} = \frac{n!}{k!(n-k)!}$, the probability of m_{C1} task control and m_{D1} data channels collided given that BF and RF task choose from previously uncollided and unjammed task channels is:

$$^{439} \quad p(m_{C1}, m_{D1} | I_C^t, I_D^t, J_C^t, J_D^t, a_B^t, a_R^t) = \frac{\binom{\alpha_{C1}^t}{m_{C1}}\binom{\alpha_{D1}^t}{m_{D1}}\binom{N_1^t - \alpha_{C1}^t - \alpha_{D1}^t}{\alpha_{I1}^t + \alpha_{J1}^t - m_{C1} - m_{D1}}}{\binom{N_1^t}{\alpha_{I1}^t + \alpha_{J1}^t}}$$

The probability of m_{C2} control and m_{D2} data channels collided given that BF and RF networks choose from *previously collided* channels is:

$$^{444} p(m_{C2}, m_{D2}|I_C^{I}, I_D^{I}, J_C^{I}, J_D^{I}, a_B^{I}, a_R^{I}) = \frac{\binom{\alpha_{C2}^{l}}{m_{C2}^{l}}\binom{\alpha_{D2}^{l}}{\alpha_{D2}^{l}}\binom{N_2^{l} - \alpha_{C2}^{l} - \alpha_{D2}^{l}}{\binom{N_2^{l}}{\alpha_{D2}^{l}}}{\binom{N_2^{l}}{\alpha_{D2}^{l}}}$$

The probability of m_{C3} control and m_{D3} data channels collided given that BF and RF networks choose from *previously* $_{448}$ *jammed* channels is:

$$^{449} p(m_{C3}, m_{D3} | I_C^t, I_D^t, J_C^t, J_D^t, a_B^t, a_R^t) = \frac{\binom{\alpha_{C3}}{m_{C3}}\binom{\alpha_{D3}}{m_{D3}}\binom{N_3^t - \alpha_{C3}^t - \alpha_{D3}^t}{\alpha_{J_2}^{t_2} - m_{C3} - m_{D3}}}{\binom{N_3^t}{\alpha_{J_2}^{t_2}}}$$

The probability of n_{C1} control and n_{D1} data channels 452 jammed given that BF and RF networks choose from 453 *previously uncollided and unjammed* channels is:

$${}^{454} p(n_{C1}, n_{D1} | I_C^t, I_D^t, J_C^t, J_D^t, a_B^t, a_R^t) = \frac{\binom{\alpha_{C1}^t}{n_{C1}}\binom{\alpha_{D1}^t}{n_{D1}}\binom{N_1^t - \alpha_{C1}^t - \alpha_{D1}^t}{\alpha_{I_1}^t + \alpha_{J_1}^t - n_{C1} - n_{D1}}}{\binom{N_1^t}{\alpha_{I_1}^t + \alpha_{J_1}^t}}$$

The probability of n_{C2} control and n_{D2} data channels 457 jammed given that BF and RF networks choose from 458 *previously collided* channels is:

$${}^{459} \quad p(n_{C2}, n_{D2} | I_C^t, I_D^t, J_C^t, J_D^t, a_B^t, a_R^t) = \frac{\binom{\alpha_{C2}^t}{n_{C2}} \binom{\alpha_{D2}^t}{n_{D2}} \binom{N_2 - \alpha_{C2} - \alpha_{D2}}{\alpha_{I_2}^t - n_{C2} - n_{D2}}}{\binom{N_2^t}{\alpha_{I_2}^t}}$$

⁴⁶⁰ The probability of n_{C3} control and n_{D3} data channels ⁴⁶¹ jammed given that BF and RF networks choose from *pre*-⁴⁶² *viously jammed* channels is:

$${}^{_{463}} p(n_{C3}, n_{D3} | I_C^t, I_D^t, J_C^t, J_D^t, a_B^t, a_R^t) = \frac{\binom{\alpha_{C3}^t}{n_{C3}}\binom{\alpha_{D3}^t}{n_{D3}}\binom{N_3^t - \alpha_{C3}^t - \alpha_{D3}^t}{\alpha_{J_2}^t - n_{C3} - n_{D3}}}{\binom{N_3^t}{\alpha_{J_2}^t}}$$

c) Posterior conditional probabilities: The combinatorial analysis leads to the posterior state transition probability distribution presented in Eq. (5), as shown at the top of next page. To solve for an optimal strategy, we need to evalutes ate this posterior distribution. For large networks and rapid variation in the system parameters (e.g., changing number of channels), this approach imposes high computational cost. ⁴⁷⁰ We can alternatively sample the distribution, using a statistically rigorous technique such as Markov Chain Monte ⁴⁷² Carlo (MCMC); however, the MCMC performance relies on ⁴⁷³ the choice of a proposal distribution that must work well for ⁴⁷⁴ CCRN, which by itself is an active area of research. In the ⁴⁷⁵ next section, we propose Q-learning [11] based methods that ⁴⁷⁶ called *value iteration* [18]. In Section V-B, we will numerically ⁴⁷⁷ evaluate transition probabilities for the two-CCRN scenario ⁴⁷⁰ using Eq. (5) and Monte Carlo sampling. ⁴⁸⁰

D. Optimal Strategies

The goal of the state-agnostic CCRN is to minimize the $_{482}$ growth of regret with an optimal strategy σ^* : $_{483}$

481

$$\sigma^* = \arg\min_{\sigma} \Gamma^t = \min_{\sigma} \left\{ \mathbb{E} \left[\sum_{i=1}^M \sum_{j=1}^t r_{(i)}^j \right] - \mathbb{E} [R_{\sigma}^t] \right\} \quad (6) \quad _{484}$$

where Γ^t represents the CCRN regret at time t, and we use $r_{(i)}^t$ ⁴⁸⁵ an *ordered* sequence of the N instantaneous channel rewards ⁴⁸⁶ at time t such that $r_{(1)}^t \ge r_{(2)}^t \ge \ldots \ge r_{(N)}^t$. Knowing that there ⁴⁸⁷ are M = C + J total number of nodes in the network, sum of ⁴⁸⁸ the M (< N) highest rewarding channels reflects the optimal ⁴⁸⁹ allocation of the nodes. We note that one can adopt the concept ⁴⁹⁰ of *discounted* reward by expressing $\mathbb{E}[R_{\sigma}^t] = \mathbb{E}\sum_{j=0}^t \gamma^j r_{\sigma}^j$. ⁴⁹¹ By adjusting the value of discount ratio γ ($0 \le \gamma < 1$) over ⁴⁹² time, one can control the coverage to which a strategy tries to ⁴⁹³ optimize short term versus long term (i.e., exploit or explore) ⁴⁹⁴ rewards.

The goal for the stateful CCRN is stated differently. We $_{496}$ wish to find an optimal distribution π^* over all possible $_{497}$ node actions to maximize the *expected* cumulative sum of $_{498}$ discounted rewards: $_{499}$

$$\pi^* = \arg\max_{\pi} \mathbb{E}\left[\sum_{j=0}^{t} \gamma^j R\left(s^j, a_B^j, a_R^j\right)\right]$$
(7) 500

where γ again is a reward discount ratio, strategy π decides BF 501 node actions, RF node actions are measurable to determine the 502 state, and the reward can be observed over time. It would be 503 an interesting problem to see whether the stateful optimization 504 of Eq. (7) is compatible to that of state-agnostic in Eq. (6). 505 We will later provide an in-depth, comparative analysis of the 506 two approaches. 507

IV. ALGORITHMIC APPROACHES TO 508 FIND OPTIMAL CCRN STRATEGIES 509

An optimal CCRN strategy yields node actions that maximize networking efficiency measurable in data throughput and jammed enemy communications. This section describes our algorithmic approaches in state-agnostic and stateful settings. We first consider the MAB formulation for its simplicity in computing strategies from only reward observation over time. In this section, we show that constructing a MAB algorithm (Thompson sampling in particular) with an extreme-valued reward likelihood as guide is ideal because the maximum possible reward should consist of only successful transmission and 519

$$p(t_{C}^{t+1}, t_{D}^{t+1}, J_{C}^{t+1}, J_{D}^{t+1} | t_{C}^{t}, I_{D}^{t}, J_{C}^{t}, J_{D}^{t}, a_{B}^{t}, a_{R}^{t})) = \sum_{\substack{l_{C}^{t+1} = m_{C1} + m_{C2} + m_{C3} \\ l_{D}^{t+1} = m_{D1} + m_{D2} + m_{D3} \\ J_{D}^{t+1} = m_{D1} + m_{D2} + m_{D3} \\ P(m_{C3}, m_{D3} | t_{C}^{t}, I_{D}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \times p(m_{C1}, n_{D1} | l_{C}^{t}, I_{D}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \\ \times p(m_{C3}, m_{D3} | t_{C}^{t}, I_{D}^{t}, J_{C}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \times p(n_{C3}, n_{D3} | l_{C}^{t}, I_{D}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \\ \times p(n_{C2}, n_{D2} | l_{C}^{t}, I_{D}^{t}, J_{C}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \times p(n_{C1}, n_{D1} | l_{C}^{t}, I_{D}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \\ \times p(m_{C3}, m_{D3} | l_{C}^{t}, I_{D}^{t}, J_{D}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \times p(n_{C3}, n_{D3} | l_{C}^{t}, I_{D}^{t}, J_{D}^{t}, d_{B}^{t}, a_{R}^{t}) \\ = \sum_{\substack{l_{C}^{t+1} = m_{C1} + m_{C2} + m_{C3} \\ J_{D}^{t+1} = m_{D1} + m_{D2} + m_{D3} \\ J_{D}^{t+1} = m_{D1} + m_{D2} + m_{D3} \\ J_{D}^{t+1} = n_{D1} + m_{D2} + m_{D3} \\ J_{D}^{t+1} = n_{D1} + m_{D2} + m_{D3} \\ \frac{l_{C}^{t} (I_{D}^{t}, I_{D}^{t}, I_{D}^{$$

⁵²⁰ jamming, not a mix of successes and failures. For comparison, ⁵²¹ we present three classical algorithms known for the stochastic ⁵²² MAB problem.

For the stateful counterpart, we apply Q-learning, a multi-agent reinforcement learning technique, in the Markov game framework. As shown previously, evaluating transition functions of a stateful CCRN is non-rivial due to combinatorial blowup. Thus, we use a value-iteration technique to evaluate the Bellman equations for our stateful algorithms. This section presents pseudo-code for all proposed algorithms.

530 A. State-Agnostic CCRN

⁵³¹ 1) Deterministic Algorithm: Algorithm 1 describes an ⁵³² adaptation of Lai and Robbins's asymptotically optimal ⁵³³ rules [7] for state-agnostic case. The algorithm keeps track of ⁵³⁴ cumulative reward $R_i^t = \sum_{j=1}^t r_i^j$ and total number of accesses ⁵³⁵ T_i^t for channel *i*. It draws two candidate channels c_{MPE} and ⁵³⁶ c_{RR} , based on the maximum point estimate (MPE) criterion ⁵³⁷ (e.g., channel with highest sample mean) and round robin ⁵³⁸ (RR) selection, respectively. The Kullback-Leibler divergence ⁵³⁹ between the two distributions serves a test statistic to finalize ⁵⁴⁰ the choice.

⁵⁴¹ 2) Indexing Algorithm: The success of Lai and Robbins ⁵⁴² depends on the accuracy of D_{KL} estimated from empirical sam-⁵⁴³ pling, which is a challenging task of its own. Another class ⁵⁴⁴ of MAB algorithms uses index as a computable substitute ⁵⁴⁵ for D_{KL} . Strictly speaking, indexing algorithm is a subclass ⁵⁴⁶ of deterministic algorithms. Auer *et al.* [19] formulated an ⁵⁴⁷ indexing scheme called Upper Confidence Bound (UCB). ⁵⁴⁸ In Algorithm 2, we present UCB for state-agnostic CCRN.

Algorithm 1 (Lai & Robbins for State-Agnostic CCRN)

1:	while $t < 1$ \triangleright initialized offline
2:	Access each channel at least once
3:	Record $R_i^t = \sum_{i=1}^t r_i^j$ and T_i^t for every channel <i>i</i>
4:	end
5:	while $t \ge 1$ \triangleright online
6:	Compute $\mu_i = R_i^t / T_i^t \ \forall i$
7:	Find MPE candidate $c_{MPE} = i^*$ s.t. $\mu_{i^*} = \max \mu_i$
8:	Find RR candidate $c_{RR} = (t \mod N) + 1$
9:	if $D_{KL}(p_{RR} \parallel p_{MPE}) > \log(t-1)/T_{C_{RR}}^t$
10:	Access c_{MPE} and observe r_{CMPE}^{t}
11:	Update R_{CMPE}^t and T_{CMPE}^t
12:	else
13:	Access c_{RR} and observe r_{CRR}^{t}
14:	Update R_{CRR}^t and T_{CRR}^t
15:	end
16:	end

Algorithm 2 (UCB for State-Agnostic CCRN)					
1: while $t < 1$	▷ initialized offline				
2: Same as Algorithm 1					
3: end					
4: while $t \ge 1$	⊳ online				
5: Compute point estimate $\mu_i = \frac{R_i^t}{T_i^t} \forall i$					
6: Compute index $g_i = \mu_i + \sqrt{\alpha \log \frac{1}{T_i^f}} \forall i$					
7: Access channel $i^* = \arg \max_i g_i$					
8: Update $R_{i^*}^t$ and $T_{i^*}^t$					
9: end					

Despite its simpler form, UCB results in the margin of error 549 that decays logarithmically in time. 550

3) Randomized Algorithm: Randomization serves an effec- ⁵⁵¹ tive means to simultaneously explore and exploit. In partic- ⁵⁵² ular, we focus on a probability matching heuristic known ⁵⁵³

Algorithm 3 (Thompson Sampling for State-Agnostic CCRN)

Require:	<i>d</i> =	=	$\{x, c$	ı, r}	for	context	<i>x</i> ,	action	а,	reward	r,	estimator
n(A d)	$) \propto$	n(rr	a f	$\frac{1}{2}n(f)$)) naran	net	erized b	NV A	9		

	$P(0 a) \circ P(0 a) \circ P(0)$	
1:	while $t \ge 1$	⊳ online
2:	Acquire x^t	
3:	Draw $\theta^t \sim p(\theta)$	
4:	Choose a^t to access $i^* = \arg \max_i \mathbb{E}[r_i^t x^t, \theta^t]$	
5:	Observe actual r^t	
6:	Update $d = d \cup \{x^t, a^t, r^t\}$	
7:	Update $p(\theta) = p(\theta d)$	

8: end

₅₅₄ as Thompson sampling [9], originally proposed in 1933 for

555 stochastic MAB problems. Thompson sampling selects actions 556 according to some optimal probability that is believed to yield ⁵⁵⁷ the maximum reward. The strategy maker based on Thompson 558 sampling is required to observe the actual outcome of a 559 selected action in order to adjust the belief. For this reason, ⁵⁶⁰ the algorithm is also called *posterior* sampling.

Despite its often superior performance over other algo-561 562 rithms, Thompson sampling lacks a rigorous theoretical 563 analysis. The full proof on the convergence properties and 564 bounds still remains to be an open problem. Thompson sam-565 pling reemerged in recent machine learning literature such 566 as Agrawal and Goyal [20], which provides the most thor-⁵⁶⁷ ough mathematical treatment available to date. Algorithm 3 describes Thompson sampling in its naturally Bayesian form. 568 We note that running Algorithms 1, 2, and 3 once deter-569 570 mines an action for one node only. For determining multiple 571 node actions, we use the following technique. Select the best channel from an algorithm and assign it to a node of interest. 572 Remove the selected channel and rerun the algorithm for the 573 574 remaining channels. Select the best channel among the remain-575 ing and assign it to the next node of interest. We repeat the 576 process until we allocate all nodes (comm or jammer).

577 B. Stateful CCRN

In the model-based reinforcement learning, an agent explic-578 579 itly learns the transition probabilities of the underlying Markov chain that characterizes the system. The model-free learn-580 ing, on the other hand, needs to evaluate the actual action 581 582 executed at a given state. We choose Q-learning [11], a model-583 free, temporal-difference learning as our baseline method to compute an optimal strategy for the stateful case. 584

Q-learning evaluates the quality of an action possible at a 585 586 particular state and the *value* of that state. The quality function 587 $Q(\cdot)$ is a function of state s and action a, and of only s for the value function $V(\cdot)$. The Bellman equations [21] characterize 589 such optimization:

590
$$Q(s, a) = R(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s')$$
(8)

$$V(s) = \max_{a'} Q(s, a')$$
⁽⁹⁾

⁵⁹² The key strength of Q-learning is the value iteration technique ⁵⁹³ that an agent performs an update $Q(s, a) = R(s, a) + \gamma V(s')$ ⁵⁹⁴ (note that s' is the next state transited from s) in place of 595 Eq. (8) without explicit knowledge of transition probability

Algorithm 4 (Q-learning for Stateful CCRN)

Require: $Q(s, a_B, a_R) = 1, V(s) = 1, \pi(s, a_B) = \frac{1}{|\mathcal{A}|}$ \forall state $s \in$ S, BF action $a_B \in \mathcal{A}$, RF action $a_R \in \mathcal{A}$; learning rate $\alpha < 1$ with decay $\lambda \leq 1$ (α , λ nonnegative) 1: while $t \ge 1$

- Draw $a_B^t \sim \pi(s^t)$ and execute 2:
- Observe r_B^t 3:
- Estimate \bar{a}_R^t given observed reward 4:
- Compute $\vec{s^{t+1}}$ 5:
- $Q(s^t, a_B^t, a_B^t) = (1 \alpha)Q(s^t, a_B^t, a_R^t) + \alpha(r_B^t + \gamma V(s^{t+1}))$ linprog: $\pi(s^t, .) = \arg \max_{a_B} \pi(s^t, a_B)Q(s^t, a_B, a_R)$ 6:
- 7:
- 8: Update $V(s^t) = \min_{a_R} \sum_{a_B} \pi(s^t, a_B) Q(s^t, a_B, a_R)$

9: Update $\alpha = \lambda \times \alpha$

10: end

p(s'|s, a), which is often too complex (e.g., Eq. (5)) to compute 596 as discussed in Section III-C. Noting that a strategy π is the 597 probability distribution of a at state s, linear programming can 598 solve for $\pi^* = \arg \max_{\pi} \sum_{a} Q(s, a) \pi$, reflecting the value 599 maximization in Eq. (9). 600

In Algorithm 4, we present a baseline Q-learning algorithm 601 that searches for an optimal strategy of the Blue Force (BF) 602 network. Note that the BF and RF networks are symmetri- 603 cal (e.g., same number of comm nodes and jammers), thus 604 have the same node action space. In Section IV-D, we pro- 605 pose three variations of Q-learning algorithms that improve 606 the performance of Algorithm 4. 607

C. New State-Agnostic Algorithm

For the state-agnostic approach, we propose a new MAB 609 algorithm based on extreme value theory [22], conjugate 610 priors, and Thompson sampling. 611

608

1) Distribution of Maximum Reward Sequence: Let $Y^t = {}_{612}$ $\max\{r_1^t, \ldots, r_N^t\}$ where r_i^t represents the reward from channel 613 *i* at *t*. Since the sequence Y^1, Y^2, \ldots, Y^t consists only of the 614 maximum channel reward each time, it must have achieved the 615 distribution p^* in Eq. (2). Furthermore, the sequence should 616 result in an *upper bound* of the optimal mean reward μ^* . 617 Therefore, we need a strategy σ to empirically follow the 618 distribution of Y^t . But how is it distributed? 619

Fisher and Tippet [23] and Gnedenko [24] proved the exis- 620 tence of limiting distributions for block maxima (or minima) 621 of random variables. Their findings became the foundation of 622 extreme value theory used widely in financial economics. 623

Theorem 1 (Fisher & Tippett, Gnedenko): Let X_1, \ldots, X_n 624 be a sequence of i.i.d. random variables and $M_n = 625$ max $\{X_1, \ldots, X_n\}$. If real number pairs (a_n, b_n) exist such that 626 $a_n, b_n > 0$ and $\lim_{n \to \infty} P(\frac{M_n - b_n}{a_n} \le x) = F(x)$, where $F(\cdot)$ is 627 a non-degenerate distribution function, then the limiting dis-628 tribution $F(\cdot)$ belongs to only Fréchet, Gumbel, or Weibull 629 family of probability distribution functions.

Proof: See Fisher and Tippet [23] and Gnedenko [24]. 631 2) Conjugate Priors: In Bayesian inference, the posterior is 632 updated by the observed likelihood given the prior distribution: 633

$$\underbrace{p(\theta|r)}_{\text{posterior}} \propto \underbrace{p(r|\theta)}_{\text{likelihood}} \times \underbrace{p(\theta)}_{\text{prior}}$$

	TABLE V	
BAYESIAN CONJUGACY	UNDER EXTREME-VA	ALUED LIKELIHOOD

Likelihood model	Conjugate priors
Pareto	Gamma
Lognormal	Gamma or normal
Exponential	Gamma
Normal	Normal
Gamma	Gamma
Weibull	Inverse gamma
Beta	Unknown
	Likelihood model Pareto Lognormal Exponential Normal Gamma Weibull Beta

Algorithm 5 (CCRN State-Agnostic Algorithm)

Require: $a_i, b_i = 0 \ \forall i$

while t < 1 ▷ initialized offline
 Access each channel until a_i, b_i ≠ 0 ∀i, where a_i and b_i are sample reward mean and variance
 end
 while t ≥ 1 ▷ online
 Draw θ_i ~ inv-gamma(a_i,b_i)

6: Estimate
$$\hat{r}_i = \text{weibull}(\theta_i, \beta_i) \forall i \text{ for given } 0.5 \le \beta_i \le 1$$

7: Access channel $i^* = \arg \max_i \hat{r}_i$

8: Observe actual r_{i*}^t to update $\{R_{i*}^t, T_{i*}^t\}$

0. Undete
$$a_{ik} = a_{ik} + T^{t}$$
 $b_{ik} = b_{ik} + \sum_{i} (r^{t})^{t}$

10: end

635 When the probabilistic model for the likelihood is known, we 636 can set the prior and posterior distributions conveniently of the 637 same family of functions. This is known as conjugate prior. 638 Since the reward distribution under our search is extremevalued, our likelihood choices are left to Fréchet, Gumbel, 639 Weibull distributions. Table V summarizes the conjugate 640 Or 641 priors having an extreme-valued likelihood distribution [25]. 3) The Algorithm: Algorithm 5 proposes a new state-642 643 agnostic approach. The algorithm performs Thompson sam-644 pling following an extreme-valued likelihood and updates the 645 posterior distribution based on its conjugate prior. However, we need to decide on which extreme value distribution is 646 suitable for CCRN. 647

Since both Fréchet and Gumbel distributions model *unbounded* random variables, we adopt a *Weibull* likelihood with the inverse gamma conjugate prior (see Table V), reasoning that the maximum reward value for CCRN should be *finite*. The lack of theoretical analysis on Thompson sampling makes it difficult to justify our design choice. In Section V, we show an empirical evidence that backs up our choice.

A Weibull distribution has finite endpoints. Its conjugate prior, the inverse gamma distribution, has two hyperparameters a, b > 0. Our algorithm draws the scale parameter θ from the inverse gamma prior $p(\theta|a, b) = \frac{b^{a-1}e^{-b\theta}}{\Gamma(a-1)\theta^a}$ for $\theta > 0$ where *a* and *b* are the sample mean and variance of the reward of a channel, and $\Gamma(\cdot)$ the gamma function (not to be confused with the Lai & Robbins's regret Γ in Eq. (1)). The Weibull random variable generated by θ drawn from the prior estimates the expected reward for the channel. After observing the actual reward, the posterior update follows.

665 D. New Stateful Algorithm

⁶⁶⁶ Before presenting the new stateful algorithm, we decom-⁶⁶⁷ pose the CCRN Markov game into two subgames, namely



Fig. 4. Inter-node relationships in antijamming and jamming subgames.

antijamming and *jamming*, and examine inter-node relationship. Fig. 4 illustrates the decomposition. Each circle represents a unique node type from the Blue Force (BF) and Red 670 Force (RF) networks. For clarity, we explain the relationships 671 in the BF network's perspective—i.e., the same relationships 672 can be derived for the RF network. 673

Antijamming game is primarily played between a BF comm 674 node and an RF jammer. That is, the key objective of anti- 675 jamming game for BF is to maximize data throughput by 676 avoiding hostile RF jamming. Antijamming game is also 677 played between a BF comm node and an RF comm node since 678 the BF comm throughput is subject to degradation when the 679 both comm nodes collide on the same channel. Fig. 4 also 680 suggests that there are antijamming games among BF comm 681 nodes and between a BF comm node and a BF jammer. When 682 multiple BF comm nodes collide on the same channel, there 683 will be a similar loss of throughput. Moreover, the BF network 684 wastes corresponding comm resources because of the colli- 685 sion, thus it should be avoided. We call misjamming when a 686 BF jammer jams BF's own comm nodes. Misjamming is dev- 687 astating because it incurs a loss of throughput and wastes both 688 comm and jamming resources of the network. 689

To avoid the collision among the BF nodes and misjamming, 690 we rely on network control to coordinate the node actions. A 691 rational strategy must first check any conflicting node actions 692 within the same network. However, imperfect sensing and 693 signaling can still lead to a collision and misjamming. 694

In jamming game, a BF jammer is trying to jam an RF 695 comm node in order to suppress the RF data throughput. A 696 BF jammer can target a data channel frequently accessed by 697 the RF comm nodes. Alternatively, it can aim for an RF control 698 channel, which would result a small immediate reward but a 699 potentially larger value in the future by blocking subsequent 700 RF data traffic. Additionally, jamming game is played between 701 a BF jammer and an RF jammer, and among multiple BF 702 jammers in order to minimize energy resources. There will be 703 no need for a BF jammer to act on a channel already under 704 jamming by either another BF jammer or an RF jammer. 705

Now we describe three comparable, new algorithms for 706 stateful CCRN using Minimax [26], Nash [27], and Friend-707 or-foe [28] Q-learning methods. 708

1) Minimax-Q Learning: Minimax-Q assumes a zero-sum 709 game that implies $Q_B(s^t, a_B^t, a_R^t) = -Q_R(s^t, a_B^t, a_R^t) = 710$ $Q(s^t, a_R^t, a_R^t)$. This holds tightly for the jamming subgame 711

$$V(s^{t}) = \max_{\pi_{B1}(A_{B,comm})} \min_{a_{R,jam}^{t}} \max_{\pi_{B2}(A_{B,jam})} \min_{a_{R,comm}^{t}} \sum_{a_{B}^{t}} Q(s^{t}, a_{B}^{t}, a_{R}^{t}) \pi_{B}(a_{B}^{t})$$
(10)

$$Q(s^{t}, a^{t}_{B}, a^{t}_{R}) = r(s^{t}, a^{t}_{B}, a^{t}_{R}) + \gamma \sum_{s^{t+1}} T(s^{t}, a^{t}_{B}, a^{t}_{R}, s^{t+1}) V(s^{t+1})$$

= $r(s^{t}, a^{t}_{B}, a^{t}_{R}) + \gamma \sum_{s^{t+1}} p(s^{t+1}|s^{t}, a^{t}_{B}, a^{t}_{R}) V(s^{t+1})$ (11)

$$Q(s^{t}, a^{t}_{B}, a^{t}_{R}) = (1 - \alpha^{t})Q(s^{t}, a^{t}_{B}, a^{t}_{R}) + \alpha^{t} \Big[r(s^{t}, a^{t}_{B}, a^{t}_{R}) + \gamma V(s^{t+1}) \Big]$$

$$Q(s^{t}, a^{t}_{B}, a^{t}_{R}) = (1 - \alpha^{t})Q(s^{t}, a^{t}_{B}, a^{t}_{R})$$
(12)

$$+ \alpha^{t} \left[r(s^{t}, a^{t}_{B}, a^{t}_{R}) + \gamma \max_{\pi_{B1}(A_{B,comm})} \min_{a^{t}_{R,jam}} \max_{\pi_{B2}(A_{B,jam})} \min_{a^{t}_{R,comm}} Q(s^{t}, a^{t}_{B}, a^{t}_{R}) \pi_{B}(a^{t}_{B}) \right]$$
(13)

$$V(s^{t}) = \max_{\pi_{B1}(A_{B,comm})} \min_{\hat{\pi}_{R2}(A_{R,jam})} \max_{\pi_{B2}(A_{B,jam})} \min_{\hat{\pi}_{R1}(A_{R,comm})} \sum_{a_{B}^{t}} Q(s^{t}, a_{B}^{t}, a_{R}^{t}) \pi_{B}(a_{B}^{t}) \hat{\pi}_{R}(a_{R}^{t}),$$
(14)

$$Q(s^{t}, a^{t}_{B}, a^{t}_{R}) = (1 - \alpha^{t}) Q(s^{t}, a^{t}_{B}, a^{t}_{R}) + \alpha^{t} \left[r(s^{t}, a^{t}_{B}, a^{t}_{R}) + \gamma \max_{\pi_{B1}(A_{B,comm}) \hat{\pi}_{R2}(A_{R,jam})} \min_{\pi_{B2}(A_{B,jam}) \hat{\pi}_{R1}(A_{R,comm})} Q(s^{t}, a^{t}_{B}, a^{t}_{R}) \pi_{B}(a^{t}_{B}) \hat{\pi}_{R}(a^{t}_{R}) \right]$$
(15)
$$Q_{B}(s^{t}, a^{t}_{B}, a^{t}_{R}) = (1 - \alpha^{t}) Q_{B}(s^{t}, a^{t}_{B}, a^{t}_{R})$$

$$+ \alpha^{t} \left[r(s^{t}, a^{t}_{B}, a^{t}_{R}) + \gamma \max_{\pi_{B1}(A_{B,comm}) \ \hat{\pi}_{R2}(A_{R,jam})} \max_{\pi_{B2}(A_{B,jam}) \ \hat{\pi}_{R1}(A_{R,comm})} \min_{\hat{\pi}_{R1}(A_{R,comm})} \mathcal{Q}_{B}(s^{t}, a^{t}_{B}, a^{t}_{R}) \ \pi_{B}(a^{t}_{B}) \ \hat{\pi}_{R}(a^{t}_{R}) \right]$$
(16)

$$\mathcal{Q}_{R}(s, a_{B}, a_{R}) = (1 - \alpha) \mathcal{Q}_{R}(s, a_{B}, a_{R}) + \gamma \max_{\pi_{B1}(A_{B,comm}) \hat{\pi}_{R2}(A_{R,jam})} \min_{\pi_{B2}(A_{B,jam}) \hat{\pi}_{R1}(A_{R,comm})} \mathcal{Q}_{R}(s^{t}, a^{t}_{B}, a^{t}_{R}) \pi_{B}(a^{t}_{B}) \hat{\pi}_{R}(a^{t}_{R}) \right]$$
(17)

The the jammer's gain is offset precisely by the opponent ming subgames jointly, we propose a slight modification to ming subgames jointly, we propose a slight modification to the original Minimax-Q algorithm in Littman [26]. First, we divide the strategy of BF network π_B into its antijamming and minimax operator to our value function in Eq. (10), as shown the top of this page. The modified Q-function in Eq. (11), as shown at the top of this page, can be computed iteratively, using Eqs. (12) and (13), as shown at the top of this page. Learning rate α decays over time such that $\alpha^{t+1} = \alpha^t \cdot \delta$ according to decay factor $0 < \delta < 1$.

2) Nash-Q Learning: Nash-Q [27] can solve a general-sum 724 725 game in addition to zero-sum games. This makes an important 726 distinction to Minimax-Q although the Nash-Q value function 727 for a zero-sum game in Eq. (14), as shown at the top of this ⁷²⁸ page, is different from Eq. (10) by only one extra term $\hat{\pi}_R(a_R^t)$. 729 This means that Nash-Q requires to estimate the policy of ⁷³⁰ the opponent. BF network needs to learn $\hat{\pi}_{R1}$ and $\hat{\pi}_{R2}$, the antijamming and jamming substrategies of RF network. The 731 732 Q-function for the zero-sum Nash-Q is given by Eq. (15), as 733 shown at the top of this page. For a general-sum game, the ⁷³⁴ BF agent should compute Q_B and Q_R separately at the same 735 time while observing its reward and estimating the RF's by 736 Eqs. (16) and (17), as shown at the top of this page. The 737 objective of Nash-Q is to find a joint equilibrium under the 738 mixed strategies $(\pi_B, \hat{\pi}_R)$.

3) Friend-or-Foe Q-Learning: Although Nash-Q is applicable to both zero-sum and general-sum games, its convergence 740 guarantee is considered too restrictive [28]. Littman [28] 741 instead proposed Friend-or-foe Q-learning (FFQ). FFQ is a 742 computational enhancement and provides better convergence 743 properties by relaxing the restrictive conditions of Nash-Q. 744 For this relaxation, FFQ requires extra information that other 745 agents in the game should be classified *friendly* or *hostile*. 746

In FFQ, the BF agent maintains only one Q-function: 747

$$Q_B(s^t, a_B^t, a_R^t) = (1 - \alpha^t) Q_B(s^t, a_B^t, a_R^t)$$

$$+ \alpha^t [r(s^t, a_B^t, a_R^t) + \gamma \Psi_B]$$
(18) 749
(18) 749

If the BF agent encounters an agent that is identified as a 750 friend, the Q-function for the BF network is updated by 751

$$\Psi_B = \max_{a_B^t, a_R^t} Q_B(s^t, a_B^t, a_R^t)$$
(19) 752

On the other hand, if the BF agent encounters an agent that 753 is identified as a foe, the Q-function is updated under the 754 minimax criterion 755

$$\Psi_B = \max_{\pi_B(A_B)} \min_{\hat{\pi}_R(A_R)} \sum_{a_B^t} Q_B(s^t, a_B^t, a_R^t) \pi_B(a_R^t).$$
(20) 756

V. SIMULATION RESULTS 757

We evaluate the state-agnostic and stateful approaches 758 against non-cognitive static and random strategies. Then, we 759



Fig. 5. Illustration of ideal channel access strategy against static RF network strategy when C = 4 and J = 2 in centralized scenario.

760 evaluate the case where one of the two CCRNs employs the ⁷⁶¹ state-agnostic algorithm and stateful algorithm for the other.

762 A. Evaluation of State-Agnostic Approach

1) Scenarios and Metric: We evaluate the centralized 763 764 and distributed control scenarios for BF and RF networks in a custom-built MATLAB simulator. The BF net-765 work runs Algorithms 1 (Lai and Robbins), 2 (UCB), 3 766 (Thompson Sampling), and 5 (proposed) described in 767 Sections IV-A and IV-C while the RF network is configured 768 with static and uniformly random strategies. In static strategy, 769 RF nodes initially choose to access some channels and con-770 tinue to access the same channels throughout. Random strategy 771 chooses a uniformly random channel for each RF node at each 772 773 time slot.

Under the centralized scenario, we assume that the central 774 775 decision maker has perfect knowledge (i.e., sensing results from all nodes) in the decision making as expressed by Eq. (3). 776 Under the distributed scenario, each node in the network 777 makes its own decision by using its sensing results only (no 778 ⁷⁷⁹ information sharing) as described by Eq. (4).

We adopt the average reward per channel as the performance 780 781 evaluation metric for a CCRN:

$$\bar{R}^t = \frac{1}{N \cdot t} \sum_{j=1}^t \sum_{i=1}^N r_i^j$$

782

⁷⁸³ where r_i is the *i*th channel reward, and there are N channels 784 in the spectrum. We use the channel reward model described by Table I in Section II-D. 785

2) Results: The spectrum has N = 10 channels. For each 786 787 CCRN, we vary the number of comm nodes C = 2, 4, 6, 8, 788 but fix the number of jammers to J = 2. Comm nodes ₇₈₉ have a transmit probability $p_{Tx} = 0.5$ whereas jammers jam 790 with probability 1. We run t = 1,000 time slots and mea-791 sure steady-state, cumulative average reward per channel for 792 comparison.

To better understand simulated results, we present an illus-793 794 trative example for an optimal BF strategy against the static 795 RF network with C = 4 and J = 2 in Fig. 5. In this example, 796 RF comm nodes are fixed at channels 1, 2, 3, 4, and its 2 797 jammers at channels 5 and 6, leaving the rest of channels 7, 798 8, 9, 10 free of RF actions. Through learning by sensing all 799 channels over time, an optimal BF strategy should place its 800 two jammers somewhere between channel 1 and 4. Because ⁸⁰¹ the comm transmit probability at any given slot is 0.5 for 802 all RF and BF comm nodes, the maximum average reward ⁸⁰³ earned by the BF jammers should be $\mathbb{E}[R_{B,jam}] \approx 0.5 \times 2 = 1$. ⁸⁰⁴ The BF comm nodes at channels 7, 8, 9, and 10 should earn





Against Static

Fig. 6. Average reward performance comparison in centralized scenario for state-agnostic approaches.

 $\mathbb{E}[R_{B,comm}] \approx 0.5 \times 4 = 2$ (because of the comm transmit 805 probability of 0.5, the BF comm reward at each time slot 806 come from two channels on the average). In summary, the total 807 reward for BF network in this example is approximately 3, 808 which is normalized to $\frac{3}{N} = \frac{3}{10} = 0.3$ (per channel), and 809 similarly for RF network, average total reward is $\frac{1}{10} = 0.1$. 810

Fig. 6 compares the performance of tested algorithms in the 811 centralized scenario. The curves for static and random strate- 812 gies were obtained while they were tested against the proposed 813 algorithm. We can clearly observe performance advantage of 814 our algorithm over Algorithms 1, 2, and 3. The proposed 815 algorithm (i.e., Algorithm 5) can learn static transmission and 816 jamming patterns effectively. Static strategy yields near-zero 817 reward at C = 2. As we have fixed J = 2, static strategy can 818 realize nonzero rewards when C > 2. 819

The results for the static and proposed algorithms at $C = 4_{820}$ in Fig. 6 confirm the reward performance of the example 821 illustrated in Fig. 5. The proposed algorithm consistently outperforms the others including the ones adapted from existing 823 MAB solutions. As the number of comm nodes per net- 824 work increases, the state-agnostic CCRN strategies face fewer 825 options in choosing node actions. As a result, we observe 826 that difference in the reward performance of the algorithms 827 becomes smaller. To increase the potential reward, we need to 828 explore more channels in optimizing the strategy. 829

Learning is harder against random strategy because ran- 830 domization gives an effective exploration mechanism. As a 831 result, the reward performance of the algorithms decreases. 832 Random strategy, however, can only explore. The lack of 833 exploitation explains its poorer performance compared to 834 Algorithms 1, 2, 3, and 5. 835

In Fig. 7, we compare the performance in the distributed 836 scenario that lacks explicit intra-network coordination. The 837 reward performance becomes worse for all of the algorithms 838 due to collisions and misjamming. 839

840

B. Evaluation of Stateful Approach

1) Scenarios and Metric: As discussed earlier, there are 841 model-based and model-free stateful approaches. We use the 842 Markov game model derived in Section III-C to numerically 843 simulate the CCRN state transition probabilities. We draw a 844 reduced state diagram in Fig. 8, using only the 10 most prob- 845 able state transitions computed from Eq. (5). Recall that a 846 CCRN state is represented by the tuple $\langle I_C, I_D, J_C, J_D \rangle$, where ⁸⁴⁷



Fig. 7. Average reward performance comparison in distributed scenario for state-agnostic approaches.



Fig. 8. Top 10 state transitions computed analytically using Eq. (5).



Fig. 9. Top 10 state transitions determined from Monte Carlo simulation (shaded states different from the analytical result of Fig. 8).

⁸⁴⁸ I_C is the number of collided control channels, I_D the number ⁸⁴⁹ of collided data channels, J_C the number of jammed control ⁸⁵⁰ channels, and J_D the number of jammed data channels.

In Fig. 9, we draw a similar state diagram, which is resulted 851 852 from Monte Carlo sampling. We notice some differences 853 between the two. First, the state diagram from the analyti-854 cal evaluation is more compact and concentrated as it has two ⁸⁵⁵ fewer states (see the two highlighted states in Fig. 9) and larger 856 transition probability values overall. One reason for larger probability values in Fig. 8 is that our Monte Carlo simulation 857 858 implements uniform sampling from the action space at each visited state, which does not necessarily reflect the behavior 859 of a rational strategy maker. Subsequently, the resulting tran-860 sitions are more likely to be distributed with less probability 861 measures. As mentioned earlier, we can improve our Monte 862 863 Carlo approach with MCMC. However, the optimal MCMC ⁸⁶⁴ design for stateful CCRN is out of scope of this paper.



Fig. 10. Average reward performances of Minimax-Q, Nash-Q, and FFQ at BF network against *static* strategy at RF network.

For model-free approach, we configure BF network to run ⁸⁶⁵ strategies based on Q-learning. We apply Minimax-Q and ⁸⁶⁶ Nash-Q learning algorithms under the centralized network ⁸⁶⁷ control while we use FFQ for the distributed case. RF network (for both model-based and model-free) is configured to ⁸⁶⁹ run the same static and random strategies that we use to evaluate the state-agnostic case. Lastly, we use the same metric as ⁸⁷¹ the state-agnostic case to evaluate the reward performance of ⁸⁷² the stateful algorithms. ⁸⁷³

The simulation parameters are as follows. There are $N = 10_{874}$ channels in the spectrum. Both BF and RF networks have $_{875}$ 2 comm nodes and 2 jammers. We set each comm node's $_{876}$ Tx and each jammer's jamming with probability 1. We sim- $_{877}$ ulate each run for 2,000 time slots and observe reward $_{878}$ performances.

2) Results: We plot the average reward performances for ⁸⁸⁰ BF network employing Minimax-Q, Nash-Q, and FFQ against ⁸⁸¹ RF network's *static* strategy over time in Fig. 10. The solid ⁸⁸² curve shows the result under the model-free approach based ⁸⁸³ on value-iterated (VI) Q-learning. The dashed curve shows the ⁸⁸⁴ result under the model-based approach using the state transitions of Fig. 8 (i.e., numerical evaluation of Eq. (5)). In the ⁸⁸⁶ steady state, we observe that the reward performances of the ⁸⁸⁷ two approaches converge. ⁸⁸⁸

The results for Minimax-Q and Nash-Q, which are obtained under the centralized network control, indicate the optimal performance against the static strategy. On the other hand, the FFQ strategy displays suboptimal reward performance due to possible collisions and misjamming in the distributed network control.

Fig. 11 presents the average reward performances for BF ⁸⁹⁵ network employing Minimax-Q, Nash-Q, and FFQ against RF ⁸⁹⁶ network's *random* strategy. The performance of Q-learning ⁸⁹⁷ algorithms decreases against the random strategy as observed ⁸⁹⁸ in the state-agnostic case. ⁸⁹⁹



Fig. 11. Average reward performances of Minimax-Q, Nash-Q, and FFQ at BF network against *random* strategy at RF network.

900 C. State-Agnostic vs. Stateful Approaches

We now evaluate the reward performance for a scenario when one network employs the state-agnostic strategy and Minimax-Q and the RF network with Algorithm 5 with simulation parameters N = 10, C = 4, and J = 2. We vary the comm Tx and jamming probabilities from 0 to 1 for the network control mechanisms.

Fig. 12 depicts the reward performances of the two networks as a function of comm Tx and jamming probabilities under the centralized control. The performances of both networks are comparable as the two approaches in the steady-state seem to achieve similar learning. In Fig. 13, we show the reward perstate formances of the two networks under the distributed control. Again, the performances are on par because both networks can end learn about each other's strategy under the distributed case.

VI. RELATED WORK

917

The state-agnostic approach of this paper is developed under 918 919 the stochastic MAB framework. In 1933, Thompson [9] intro-920 duced a stochastic MAB problem and proposed an optimal 921 heuristic known as Thompson sampling, which remains to ⁹²² be an effective action selection strategy. Robbins [29] 1952 923 presented the first sequential analysis of the single-player 924 MAB problem. In Bellman [30] 1954, MAB problems were 925 formulated as a class of Markov decision process (MDP). 926 Gittins [31] 1979 proved the existence of a Bayes opti-927 mal indexing scheme for MAB problems if they can be 928 modeled as a stationary MDP. Lai and Robbins [7] 1985 introduced the notion of regret, derived its lower bound using the 929 930 Kullback-Leibler divergence, and constructed asymptotically ⁹³¹ optimal allocation rules. Anantharam et al. [16] 1987 extended 932 Lai & Robbins for multi-player. Whittle [32] 1988 intro-933 duced PSPACE-hard restless MAB problems and showed that



Fig. 12. Average reward performances of Minimax-Q (stateful approach) at BF network against Algorithm 4 (state-agnostic approach) at RF network under the centralized network control.



Fig. 13. Average reward performances of Minimax-Q (stateful approach) at BF network against Algorithm 4 (state-agnostic approach) at RF network under the distributed network control.

suboptimal indexing schemes are possible. Rivest and Yin [33] ⁹³⁴ 1994 proposed Z-heuristic that achieved a better empirical ⁹³⁵ performance than Lai and Robbins. Auer *et al.* [19] 2002 proposed Upper Confidence Bound (UCB), an optimistic indexing ⁹³⁷ scheme. ⁹³⁸

The foundation of our stateful approach is reinforcement ⁹³⁹ learning [34], which extends beyond self-confined views of the ⁹⁴⁰ classical Markov Decision Process in which an agent's environment is stationary and contains no other agents. Q-learning ⁹⁴² ⁹⁴³ was originally proposed by Watkins and Dayan [11]. 944 Littman [26] proposed Minimax-Q learning for a zero-sum 945 two-player game. Littman and Szepesvári [18] showed that Minimax-Q converges to the optimal value suggested by game 946 947 theory. Hu and Wellman [27] described Nash-Q that was dis-⁹⁴⁸ tinguished from Minimax-Q by solving a general-sum game with a Nash equilibrium computation in its learning algorithm. 949 Nash-Q has more general applicability, but its assumptions 950 on the sufficient conditions for convergence guarantee are 951 952 known to be restrictive. Friend-or-foe Q-learning (FFQ) [28] 953 converges precisely to the steady-state value that Nash-Q 954 guarantees. The key improvement of FFQ is relaxation of 955 the restrictive conditions that Nash-Q has, but FFQ requires priori knowledge on other agents identified as either a 956 a 957 friend or foe.

This paper considers some similar problems discussed by 958 Wang et al. [2] such as finding a strategy against hostile 959 960 jamming. They formulated a stochastic antijamming game ⁹⁶¹ played between the secondary user and a malicious jammer, provided sound analytical models, and applied unmodified 962 Minimax-Q learning to solve for the optimal antijamming 963 strategy. Our work is novel and differentiated from existing 964 work by the following. We embrace the notion of friendly jam-965 mers and provide an integrated antijamming-jamming strategy 966 967 for cognitive radio networks under competition in a hostile 968 environment. We use jamming as a means to cope with crit-⁹⁶⁹ ical situations assumed in tactical mobile networking. At the 970 same time, we try to avoid the hostile jammers that pose a 971 serious threat to the network's comm activities in another opti-⁹⁷² mization. We promote the notion of strategic jamming enabled 973 by reinforcement learning. We modify existing Q-learning 974 algorithms to solve for optimal antijamming and jamming 975 strategies jointly. Lastly, we remark that this paper unites our 976 previous approaches [12], [13] for competing cognitive radio 977 networks.

978

VII. CONCLUSION

⁹⁷⁹ We have described Competing Cognitive Resilient ⁹⁸⁰ Networks (CCRNs) that operate in a hostile environment to ⁹⁸¹ maximize data throughput while simultaneously suppressing ⁹⁸² opponent's comm activities. We have provided two compati-⁹⁸³ ble, yet distinguished approaches for an optimal strategy that ⁹⁸⁴ can coordinate the joint comm and jamming actions for a ⁹⁸⁵ CCRN.

In our state-agnostic approach, we have adopted the MAB framework and extended classical solutions for CCRN. In addition, we have proposed a new algorithm that outperforms between the classical solutions. The new algorithm builds on Thompson sampling in its Bayesian form with an enhancement from the extreme value theory. Our performance results indicate that the proposed state-agnostic algorithm proves to be most effecsite in addressing the exploration-exploitation tradeoff when even to other algorithms.

For the stateful approach, we have modeled the dynamics of CCRNs using the Markov game framework and decomposed it to antijamming and jamming subgames. We have derived an analytical expression to evaluate model-based reinforcement learning that can solve for an optimal strategy. ⁹⁹⁹ We also have applied model-free reinforcement learning, ¹⁰⁰⁰ namely Minimax-Q, Nash-Q, Friend-or-foe Q (FFQ) algo- ¹⁰⁰¹ rithms, to achieve optimal strategies. With Monte Carlo ¹⁰⁰² sampling, we have demonstrated that the model-based and ¹⁰⁰³ model-free methods approach to the same steady-state value. ¹⁰⁰⁴

Both state-agnostic and stateful approaches can be applied 1005 to centralized or distributed network control mechanisms. The 1006 numerical results suggest the superior performance achieved 1007 for the case of centralized where CCRN nodes are cooperative 1008 and their actions are coordinated through a single entity. On 1009 the other hand, the performance under the distributed control 1010 mechanism suffers from collisions and misjamming. 1011

Based on the reward performance, both state-agnostic and 1012 stateful approaches are significantly better than rudimen- 1013 tary strategies such as static and random. However, when 1014 competing against each other, the state-agnostic and stateful 1015 approaches achieve the same performance. Our future work 1016 includes various improvements for cognitive learning such 1017 as faster learning, better convergence properties, and lower 1018 computational complexity.

ACKNOWLEDGMENT

Any opinions, findings, conclusions or recommendations 1021 expressed in this material are those of the authors and do not 1022 necessarily reflect the views of the Office of the Secretary of 1023 Defense or the United States Government.

References

- Q. Zhao and B. M. Sadler, "A survey of dynamic spectrum access," 1026 IEEE Signal Process. Mag., vol. 24, no. 3, pp. 79–89, May 2007. 1027
- B. Wang, Y. Wu, K. J. R. Liu, and T. C. Clancy, "An anti-jamming 1028 stochastic game for cognitive radio networks," *IEEE J. Sel. Areas* 1029 *Commun.*, vol. 29, no. 4, pp. 877–889, Apr. 2011.
- [3] R. Chen, J.-M. Park, and J. H. Reed, "Defense against primary user emu- 1031 lation attacks in cognitive radio networks," *IEEE J. Sel. Areas Commun.*, 1032 vol. 26, no. 1, pp. 25–37, Jan. 2008.
- [4] M. Li, I. Koutsopoulos, and R. Poovendran, "Optimal jamming attack 1034 strategies and network defense policies in wireless sensor networks," 1035 *IEEE Trans. Mobile Comput.*, vol. 9, no. 8, pp. 1119–1133, Aug. 2010. 1036
- [5] S. Khattab, D. Mosse, and R. Melhem, "Modeling of the 1037 channel-hopping anti-jamming defense in multi-radio wireless net- 1038 works," in *Proc. Int. Conf. Mobile Ubiquitous Syst. Comput. Netw.* 1039 *Services (Mobiquitous)*, Dublin, Ireland, 2008, Art. no. 25. 1040
- [6] V. Navda, A. Bohra, S. Ganguly, and D. Rubenstein, "Using channel 1041 hopping to increase 802.11 resilience to jamming attacks," in *Proc. IEEE* 1042 *Int. Conf. Comput. Commun. (INFOCOM)*, Anchorage, AK, USA, 2007, 1043 pp. 2526–2530. 1044
- [7] T. L. Lai and H. Robbins, "Asymptotically efficient adaptive allocation 1045 rules," Adv. Appl. Math., vol. 6, no. 1, pp. 4–22, 1985.
- [8] L. Lai, H. El-Gamal, H. Jiang, and H. V. Poor, "Cognitive medium 1047 access: Exploration, exploitation, and competition," *IEEE Trans. Mobile* 1048 *Comput.*, vol. 10, no. 2, pp. 239–253, Feb. 2011. 1049
- W. R. Thompson, "On the likelihood that one unknown probability 1050 exceeds another in view of the evidence of two samples," *Biometrika*, 1051 vol. 25, nos. 3–4, pp. 285–294, 1933.
- [10] L. S. Shapley, "Stochastic games," Proc. Nat. Acad. Sci. USA, vol. 39, 1053 no. 10, pp. 1095–1100, 1953.
- [11] C. Watkins and P. Dayan, "Q-learning," Mach. Learn., vol. 8, no. 3, 1055 pp. 272–292, 1992.
- [12] Y. Gwon, S. Dastangoo, C. Fossa, and H. T. Kung, "Competing 1057 mobile network game: Embracing antijamming and jamming strate- 1058 gies with reinforcement learning," in *Proc. IEEE Conf. Commun. Netw.* 1059 *Security (CNS)*, 2013, pp. 28–36. 1060

1025

1020

- 1061 [13] Y. Gwon, S. Dastangoo, and H. T. Kung, "Optimizing media access
- strategy for competing cognitive radio networks," in *Proc. IEEE Glob. Commun. Conf. (GLOBECOM)*, Atlanta, GA, USA, 2013,
 pp. 1215–1220.
- W. Xu, W. Trappe, Y. Zhang, and T. Wood, "The feasibility of launching and detecting jamming attacks in wireless networks," in *Proc. ACM Int. Symp. Mobile Ad Hoc Netw. Comput. (MobiHoc)*, Chicago, IL, USA, 2005, pp. 46–57.
- 2005, pp. 46–57.
 1069 [15] M. Pajic and R. Mangharam, "Anti-jamming for embedded wireless networks," in *Proc. Int. Conf. Inf. Process. Sensor Netw. (IPSN)*, San Francisco, CA, USA, 2009, pp. 301–312.
- 1072 [16] V. Anantharam, P. Varaiya, and J. Walrand, "Asymptotically efficient allocation rules for the multiarmed bandit problem with multiple plays-
- 1074
 Part I: I.I.D. rewards," IEEE Trans. Autom. Control, vol. 32, no. 11,

 1075
 pp. 968–976, Nov. 1987.

 1076
 [17]
 T. M. Cover and J. A. Thomas, Elements of Information Theory.
- 1077 Hoboken, NJ, USA: Wiley, 1991.
- 1078 [18] M. L. Littman and C. Szepesvári, "A generalized reinforcement-learning model: Convergence and applications," in *Proc. Int. Conf. Mach. Learn. (ICML)*, Bari, Italy, 1996, pp. 310–318.
- 1081 [19] P. Auer, N. Cesa-Bianchi, and P. Fischer, "Finite-time analysis of the multiarmed bandit problem," *Mach. Learn.*, vol. 47, nos. 2–3, pp. 235–256, May/Jun. 2002.
- 1084 [20] S. Agrawal and N. Goyal, "Analysis of Thompson sampling for the multi-armed bandit problem," in *Proc. Conf. Learn. Theory (COLT)*, 2012, Art. no. 39.
- 1087 [21] R. Bellman, *Dynamic Programming*. Princeton, NJ, USA: Princeton Univ. Press, 1957.
- L. de Haan and A. Ferreira, *Extreme Value Theory: An Introduction* (Springer Series in Operations Research and Financial Engineering).
 New York, NY, USA: Springer-Verlag, 2006.
- R. A. Fisher and L. H. C. Tippett, "Limiting forms of the frequency distribution of the largest or smallest member of a sample," *Math. Proc. Camb. Philosoph. Soc.*, vol. 24, no. 2, pp. 180–190, 1928.
- B. V. Gnedenko, "Sur la distribution limite Du terme maximum D'une serie aleatoire," *Ann. Math.*, vol. 44, no. 3, pp. 423–453, 1943.
- 1097 [25] E. I. George, U. E. Makov, and A. F. M. Smith, "Conjugate likelihood
- distributions," *Scandinavian J. Stat.*, vol. 20, no. 2, pp. 147–156, 1993. 1099 [26] M. L. Littman, "Markov games as a framework for multi-agent
- reinforcement learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, New Brunswick, NJ, USA, 1994, pp. 157–163.
- 1102 [27] J. Hu and M. P. Wellman, "Multiagent reinforcement learning: Theoretical framework and an algorithm," in *Proc. Int. Conf. Mach. Learn. (ICML)*, Madison, WI, USA, 1998, pp. 242–250.
- 1105 [28] M. L. Littman, "Friend-or-foe Q-learning in general-sum games," in
- Proc. Int. Conf. Mach. Learn. (ICML), Williamstown, MA, USA, 2001, pp. 322–328.
 [29] H. Robbins, "Some aspects of the sequential design of experiments,"
- 1108 [29] H. Robbins, "Some aspects of the sequential design of experiments," Bull. Amer. Math. Soc., vol. 58, no. 5, pp. 527–535, 1952.
- 1110 [30] R. E. Bellman, A Problem in the Sequential Design of Experiments. Fort
 1111 Belvoir, VA, USA: Defense Tech. Inf. Center, 1954.
- 1112 [31] J. C. Gittins, "Bandit processes and dynamic allocation indices," *J. Roy.* 1113 *Stat. Soc.*, vol. 41, no. 2, pp. 148–177, 1979.
- 1114 [32] P. Whittle, "Restless bandits: Activity allocation in a changing world,"
 1115 J. Appl. Probab., vol. 25, pp. 287–298, Jan. 1988.
- 1116 [33] R. L. Rivest and Y. Yin, "Simulation results for a new two-armed bandit heuristic," in *Proc. Workshop Comput. Learn. Theory Nat. Learn. Syst.*,
 1118 Cambridge, U.K., 1994, pp. 477–486.
- [34] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*.
 Cambridge, MA, USA: MIT Press, 1998.



Siamak Dastangoo (S'88–M'00) received the bachelor's degree and the master's degree in electrical engineering from the University of Massachusetts, and the doctoral degree in electrical engineering from George Washington University. He spent several years in the commercial and defense industries working on a range of problems in communications and networking systems. In 2006, he joined Lincoln Laboratory, Massachusetts Institute of Technology, where he is a Technical Staff Member with the Tactical Networks Group, conducting research and

¹¹³² development in the area of wireless ad hoc networks. His current research ¹¹³³ interests are in the areas of cognitive networks and performance predictions ¹¹³⁴ of networks.



Carl E. Fossa (S'98–M'02–SM'05) received the 1135 bachelor's degree from United States Military 1136 Academy, the master's degree from the Air Force 1137 Institute of Technology, the doctoral degree from 1138 Virginia Polytechnic Institute and State University, 1139 all in electrical engineering. He served as an Army 1140 Signal Officer for 21 years, retiring at the rank 1141 of Lieutenant Colonel. In 2008, he joined Lincoln 1142 Laboratory, Massachusetts Institute of Technology, 1143 where he is an Associate Leader of the Tactical 1144 Networks Group, focussing on the performance of 1145

mobile ad hoc networks in tactical military environments. Applications of 1146 this research include the development and deployment of a real-time emula- 1147 tion of the Army Warrior Information Network-Tactical to Aberdeen Proving 1148 Grounds, MD, USA. He served in a range of tactical military positions, which 1149 included deployment to Operation Desert Shield/Storm. He also served in a 1150 number of technical engineering positions at major command headquarters 1151 and as an Assistant Professor of electrical engineering with the United States 1152 Military Academy.



Youngjune L. Gwon (S'12–M'14) received the 1154 B.S. degree from Northwestern University and the 1155 M.S. degree from Stanford University, both in electrical engineering, and the Ph.D. degree in com- 1157 puter science from Harvard University, in 2015. 1158 He is a Technical Staff Member with the Human 1159 Language Technology Group, Lincoln Laboratory, 1160 Massachusetts Institute of Technology, where he 1161 does machine learning research and development for 1162 speech and audio-video processing systems. He cur- 1163 rently works on multimodal representation learning 1164

for large, unstructured audio, video, and text corpora. His broader research 1165 interests include networks, wireless communications, embedded systems, and 1166 security. He was a Software Engineer with Silicon Valley for ten years. 1167



Hsiang-Tsung Kung received the Ph.D. degree 1168 from Carnegie Mellon University. He taught at 1169 Carnegie Mellon University for 19 years. In 1992, 1170 he joined Harvard University, where he is the 1171 William H. Gates Professor of Computer Science 1172 and Electrical Engineering. His broad research areas 1173 include complexity theory, database systems, very 1174 large scale integration, parallel computing, computer 1175 networks, security, and wireless communications. 1176 He currently focuses on machine learning, high- 1177 performance computing, and the Internet of Things. 1178

He is well known for his pioneering work on systolic arrays in parallel pro- 1179 cessing and optimistic concurrency control in database systems. His academic 1180 honors include a membership in the National Academy of Engineering and 1181 the ACM SIGOPS 2015 Hall of Fame Award (with J. Robinson) that rec- 1182 ognizes the most influential operating systems papers that were published at 1183 least ten years in the past. 1184