

March 2, 2016

## 1 Single SW

We first extend CO2 to the case with single switch and multiple receivers. The topology of the switch and receivers is shown in Figure . The switch holds rules with multiple action tag, and each action represents forwarding the traffic to the receiver with the same action tag. The receivers feedback the top overselected IPs with regard to bytes volume to controller, and the controller decides which one to put into the Blacklist in the switch. The challenge is that the Blacklist memory is a bottleneck. We propose a Reinforcement Learning (RL) based Approach (RLA) in the controller to allocate Blacklist space to the rules with different actions.

Assume the system contains 1 switch and  $n$  receivers, and we define a utility function  $U_i$  for receiver  $r$ , we aim to maximize the following objective:

$$\max \sum_{r=1}^n U_r(a_{r,t}) \text{ st. }, \\ \sum_{r=1}^n a_{r,t} \leq A_c$$

Where  $a_{r,t}$  is the Blacklist space allocated to the rules from receiver  $r$  at time  $t$ , and  $A_c$  is the Blacklist space capacity.

We use reinforcement learning to compute the utility function for each receiver, so,

$$U_r(a_{r,t}) = Q_r(s_{r,t}, a_{r,t})$$

where function  $Q_r$  is the value function in the reinforcement learning system, and  $s_{r,t}$  is the state at time  $t$ .

Reinforcement learning is learning what to do-how to map situations to actions-so as to maximize a numerical reward signal[Reinforcement book]. Beyond the agent and the environment, one can identify four main subelements of a reinforcement learning system: a policy, a reward function, a value function, and optionally, a model of the environment[Reinforcement book]. A policy defines the learning agent's way of behaving at a given time. A state  $s_{r,t}$  is the Blacklist space already allocated to the rules from receiver  $r$ , and a policy  $a_{r,t}$  is to increase or decrease the allocation given on the current state. A reward function  $r_t$  defines the goal in a reinforcement learning problem. We set a target overselection rate  $ov_T$ , and if the actual overselection rate  $ov_a$  is larger than  $ov_T$ , the reward function is a negative const value, otherwise, it is a positive const value. By performing the reward function, we aim to control the actual overselection rate under the target value.

We use Q-learning in our design, and the one-step Q-learning is defined by,

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \{r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)\}$$

## 2 Multiple SW

We extend our design to a network-wide case, with multiple switches and multiple receivers. Assume the system contains  $m$  switches, we define the objective function by,

$$\max \sum_{s=1}^m \sum_{r=1}^n U_{s,r}(a_{s,r,t}) \text{ st.},$$

$$\sum_{r=1}^n a_{s,r,t} \leq A_{c,s}$$

$$\text{And, } U_{s,r}(a_{s,r,t}) = Q_r(s_{s,r,t}, a_{s,r,t})$$

Where the subscript  $s$  indicates switch  $s$