

# Smart Resource Allocation for Mobile Edge Computing: A Deep Reinforcement Learning Approach

Pavitra Basappa Gudimani

Summer term 2020

## Abstract

The amount of data generated by all the mobile devices in the future can cause congestion at the cloud for data processing and cause a delay in service time. Mobile Edge Computing (MEC) architecture extension of existing cloud computing is promising solution. MEC improves the response time by moving network and computing resources closer to the customer. MEC does not send the data to the cloud for further processing rather the edge network itself will process the data. The burst request at MEC faces an issue in allocating computing and network resources. At this stage, Deep Reinforcement Learning based Resource Allocation (DRLRA) scheme adapts to a mutative nature and allocates computing and network resources. Thus, it reduces the service time to users and offers fairness in the use of resources at MEC. DRLRA achieved better results in varying MEC environments than the traditional OSPF algorithm.

## 1 Introduction

The phenomena of embedding the sensors in smart devices to sense the environment have changed the way of living life. The data sent by IoT devices from smart cities, smart healthcare, smart vehicles [1], etc. results in a high volume of data. And it is produced at a high pace (the European Commission has predicted that there will be 50 to 100 billion smart devices connected to the Internet by 2020) [2]. Also, Cisco has predicted that the devices connected to

the Internet will generate 507.5 ZB/year by 2019 [2]. The generated data is further transferred to the remote cloud. The volume and speed of the data generated create congestion at the cloud. Thus, efficient data processing architecture is required.

At this stage, Mobile Edge Computing architecture aid in reducing the network congestion and response delay. It deploys the applications in a distributed form in mobile edge servers to deliver the computational and critical applications services. It appears to be an encouraging solution to the problem faced by several mobile devices. On the other hand, it has following limitations. Primarily, the initial cost of investment to deploy MEC servers (MECS) at the base stations and for their maintenance is expensive. Secondly, Only a restricted number of applications can be deployed to supply services to innumerable edge devices [2]. Lastly, the dynamic nature of human life brings new requirements to the devices, continuously changing the location of requests generation, and inconsistent quantity of data, all these factors cause serious burst over MECS, increases the computational load, and congestion in the network links of certain regions.

MEC fails to provide flexibility in granting computing and network resources. Therefore, the author has presented Deep Reinforcement Learning based Resource Allocation algorithm to assign computing and network resources accordingly under the dynamic circumstances of MEC. DRLRA uses DRL technology to understand the environment and improve its behavior to perform favorable action to allocate resources. Software Defined Network with

MEC architecture centralizes the distributed network. The main goal is to minimize service time by utilizing DRLRA.

## 2 SDN-enabled MEC architecture

Deploying all the applications on a single MEC Server is expensive. MECS solves the problem by virtualizing the resources. When mobile requests for an application, it will be sent to MECS in its vicinity. If there is no required application, the request will be routed to other MECSs with the needed application deployed. Hence, the total service time includes the access time to nearby MECS, routing time among MECSs, and processing time at corresponding MECS.

To have the request delivered to destined MECS in the mutative environment the routing should happen correctly. That's where the emerging SDN technology comes in handy to understand the state of the network. In this context SDN uses below planes [3].

1. Data Plane - Forwarding devices Switches, routers, and gateways reside in this plane and work according to decisions made by the controller in the control plane.
2. Control Plane - The core of SDN architecture and all forwarding decisions made upon the control logic provided here.

DRL unit with the DRLRA algorithm is installed in the SDN controller plane and it takes dynamic states of the network as input from the data plane and it trains the agent to make decisions, based on which SDN controller can perform forwarding.

## 3 Problem Definition

MEC without resource allocation algorithm faces below problems

**Edge network routing delay:** The routing delay for the aggregated requests of mobile devices at nearby MECS to the MECS that is hosting the required application. [4].

**Data Processing delay:** Number of CPU Cycles required for computation at the destined MECS.[4].

**Variance of network resource allocation:** Several

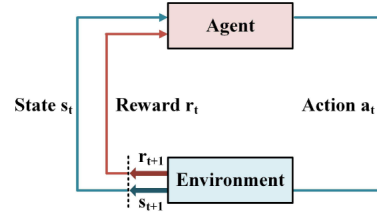


Figure 1: Illustration of Reinforcement Learning[4]

requests will be forwarded to destined MECS. Allocation of network resources on the data links should be stabilized to reduce network traffic and to transfer the request.

## 4 Deep Reinforcement Learning Based Resource Allocation Algorithm

A smart resource algorithm understands the varying environment of the MECS and makes a feasible sequence of decision on resource allocation using DQN framework.

### 4.1 Reinforcement Learning

It is all about training an agent that communicates with the environment. The agent shifts between different scenarios in the environment indicated as states, by performing actions. The action in return helps the agent to grab positive, negative or zero rewards. The sole purpose of an agent is to maximize the total rewards over the episode. The agent is reinforced to perform certain actions by giving positive rewards and keeping them away from certain actions by giving a negative reward. This way agent learn to strengthen its behavior [5]. In the given situation, the agent takes the MEC environment as a state, performs network, and computing resource allocation action. If the agent's action managed to reduce service time and balance resource allocation it receives positive rewards and vice versa [4]. Fig.1 explains, an episode at time  $t$ , the agent will get the knowledge on environment i.e state  $s_t$ , make action  $a_t$  based on policy  $\pi$  and gain best reward  $r_t$ . Further, the agent

enters new state  $s_{t+1}$ . The policy will be updated based on the reward is gained [4].

- **Episode (k):** Sequence of states, actions, and rewards that occurs between initial state and termination state. In this context, it is a resource allocation process from MECS where requests from mobile devices gathered to get a service of application to destined MECS where the required application is installed.
- **Policy  $\pi$ :** is mapping from states of the environment to actions to be taken in those states [6].
- **Cumulative Reward:  $R_t$ ,** Sum of reward obtained by an agent at the end of each episode.  $R_t = \sum_{i=t}^T \gamma^{i-t} r_i$ ,  $\gamma \in [0, 1]$ ,  $\gamma$  is discount factor to indicate the importance of future reward to the current state.  $T$  is the final step of an episode [4].
- **Action value function (Q function):** The value function is to quantify the quality of the action made at the state.

## 4.2 Deep Q Network

DQN is a powerful framework of Deep Reinforcement Learning (DRL). It takes the current state as an input and gives action value for every action that can be taken in that state i.e  $Q(s, a; \omega)$ , where  $\omega$  is a parameter of network. Action with the highest Q value  $a = \operatorname{argmax}_a Q(s, a; \omega)$  is selected using  $\epsilon$ -greedy strategy. To make the generated Q value closer to the target Q value, Mean Square Error, a loss optimizing function is used [4].

$$L(\omega) = E[(r + \gamma \max_{a_{t+1}} Q'(s_{t+1}, a_{t+1}; \omega) - Q(s_t, a_t; \omega))^2], \quad (1)$$

where  $r + \gamma \max_{a_{t+1}} Q'(s_{t+1}, a_{t+1}; \omega)$  is the target Q value and  $Q(s_t, a_t; \omega)$  is evaluation network Q value.

## 4.3 Training

DRL agent in a MEC environment is trained using DQN framework. Figure 2 explains the training process. Experience Replay memory is installed to store environment response for taken action and parameters  $\omega$  and  $\omega^-$  are set for target and evaluation network respectively. Episode k is initialized

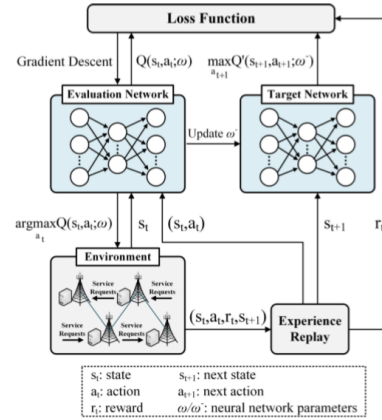


Figure 2: Detailed DQN Framework [4]

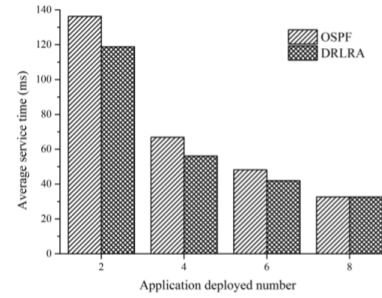


Figure 3: OSPF and DRLRA comparison under different application deployed number [4]

with state  $s_t$  at time  $t$ .  $s_t$  is taken as input to evaluation network and using the  $\epsilon$ -greedy method action  $a_t$  with maximum Q value is chosen. The agent gets a reward  $r_t$  and next state  $s_{t+1}$ .  $Q(s_t, a_t, r_t, s_{t+1})$  is stored in experience replay. To train the network, sample state  $s_t$  and action  $a_t$  applied to evaluation network from memory.  $Q(s_j, a_j; \omega)$  and  $\gamma \max_{a_{j+1}} Q'(s_{j+1}, a_{j+1}; \omega)$  calculated from evaluation and target network respectively. The loss function is used to bring the Q value of evaluation network closer to the target network through gradient descent and backpropagation. Finally, the trained evaluation network is obtained after performing resource allocation of  $t$  steps in episode  $k$  [4].

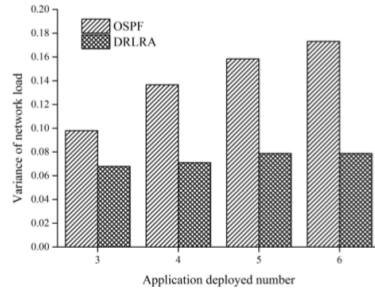


Figure 4: OSPF and DRLRA comparison under variance of network load [4]

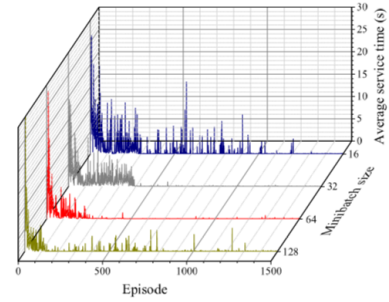


Figure 5: Convergence performance of DRLRA [4]

## 5 Performance Evaluation

**Average service time minimization:** Fig.3 shows that the average service time decrease with an increase in the number of deployed applications. As the number of deployed application increase, each MECS will almost have the requested application. Therefore, the distance between the location of requests generated for a specific application and MECS hosting the required application will reduce and thus result in a reduction in average service time. At a certain point, when the number of deployed applications reaches a maximum, the service time is considerably decreased.

**Resource Allocation Balancing Analysis:** Fig. 4 depicts that OSPF performance degrades over the increase in the complexity of the MECS environment by rising number of deployed applications. On the other hand, DRLRA, while making decisions on path selection for routing it learns the situation and adapts to it. Therefore, DRLRA allocates resources by using it evenly and prevents network congestion.

**Convergence Performance Analysis:** Fig.5 shows as the batch of samples provided increases, the speed required to stable the average service time is reduced. The agent will not get a chance to learn enough with the small number of samples. With the huge batch size agent might converge to the poor average service time. Therefore, the size of samples should be chosen carefully.

## References

- [1] Wen Sun, JiaJia Liu, and Haibin ZHANG. *When Smart Wearables meet Intelligent Vehicles: challenges and Future direction*. IEEE Wireless Communications, 2017.
- [2] Xiang Sun and Nirwan Ansari. *EdgeIoT: Mobile Edge Computing for the Internet of Things*. IEEE Communications Magazine, 2016.
- [3] Kuljeet Kaur et al. *Edge Computing in the Industrial Internet of Things Environment: Software-Defined Networks-Based Edge-Cloud Interplay*. IEEE Communications Magazine, 2018.
- [4] Jiadai Wang et al. *Smart Resource Allocation for Mobile Edge Computing: A Deep Reinforcement Learning Approach*. IEEE Transactions on Emerging Topics in Computing, 2019.
- [5] *Deep Q Networks*. URL: <https://towardsdatascience.com/q-learning-from-the-ground-up-1bbda41d3677>.
- [6] *Policy*. URL: <https://stackoverflow.com/questions/46260775/what-is-a-policy-in-reinforcement-learning>.