



JOINT INSTITUTE
交大密西根学院

DeepReserve: dynamic edge server reservation for connected vehicles with deep reinforcement learning

Jiawei Zhang

Wireless Networking and Artificial Intelligence Lab

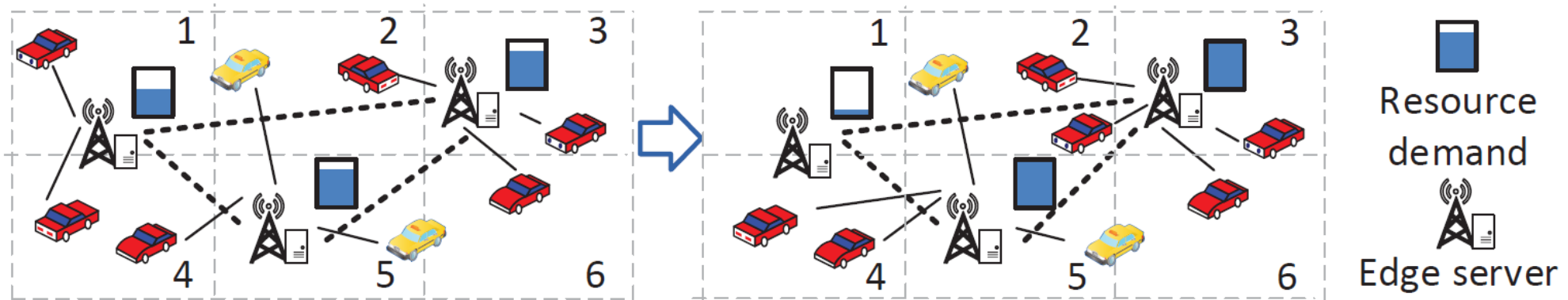
Advisor: Prof. Xudong Wang

<http://wanglab.sjtu.edu.cn>



■ Problem statement

- Computational resource demands change according to vehicle mobility



- Edge servers should be dynamically reserved (rent or start) for providing service
 - Reduce cost for idle server
 - Ensure available server for nearby users

■ Challenges

- Reservation according to statistical demands is infeasible
 - Participants are **unwilling** to share routes
 - Statistical information **cannot be aligned** with real-time demands



■ Physical placement of edge servers

■ With given demand information

- Determine the positions of edge servers that are geographically close to users [1,2]
- Determine locations of edge servers based on the number of user requests aggregated in nearby BSs [3]
- Choose suitable edge servers to hold multiple interrelated services [4]
- Determine the placement of multiple services into edge servers with heterogeneous capacities [5]

Difference: they determine **fixed placement** based on **given statistical demands**, while the edge-server reservation problem requires an **online solution** adaptive to the real-time demands

■ Without demand information

- Leverages the collected contexts of connected users (e.g., equipment types and external environment factors) to predict the demands [6]

Difference: we do not assume any context information of users

[1] H. Yin, et. al, "Edge provisioning with flexible server placement," *IEEE Trans. Parallel Distrib. Syst (TPDS)*, 2016.

[2] Z. Xu , et. al, "Efficient algorithms for capacitated cloudlet placements," *IEEE Trans. Parallel Distrib. Syst (TPDS)*, 2015.

[3] S. Wang , et. al, "Edge server placement in mobile edge computing," *J. Parallel Distrib. Computing*, 2019.

[4] I. Lera , et. al, "Availability-aware service placement policy in fog computing based on graph partitions," *IEEE Internet of Things J*, 2018.

[5] S. Pasteris , et. al, "Service placement with provable guarantees in heterogeneous edge computing systems," in *Proc. INFOCOM*, 2019.

[6] L. Chen , et. al, "Spatio-temporal edge service placement: A bandit learning approach," *IEEE Trans. Wireless Commun. (TWC)*, 2018.



■ Optimization problem

- **Task:** dynamically reserve edge servers $x_{i,t}$
- **Target:** maximize the system utility
 - (Reward of providing service) – (server cost) – (punishment of connection failure)

$$\max \sum_{i \in \mathcal{E}} (-\alpha x_{i,t} + \beta u_{i,t} - \gamma q_{i,t})$$

■ Constraints

- Latency $d_{i,j,t} y_{i,j,t} \leq D x_{i,t}$
- Connect to at most one server $\sum y_{i,j,t} \leq 1$
- Edge server capacity $\sum_{j \in \mathcal{V}_t} y_{i,j,t} \leq U x_{i,t}$
- Connect to reserved server $y_{i,j,t} \leq x_{i,t}$

■ Infeasibility of optimization method

- This problem can be reduced to the **K-median** problem which is proved to be **NP-hard** [1]
- Lack of parameters
 - Workload $u_{i,t}$ cannot be obtained before reservation
 - Latency $d_{i,j,t}$ cannot be measured, due to unknown CV location and network status

[1] M. Charikar, S. Guha, E. Tardos, and D. B. Shmoys, "A constant-factor approximation algorithm for the k-median problem," *J. Comput. And Syst. Sciences*, vol. 65, no. 1, pp. 129–149, 2002.

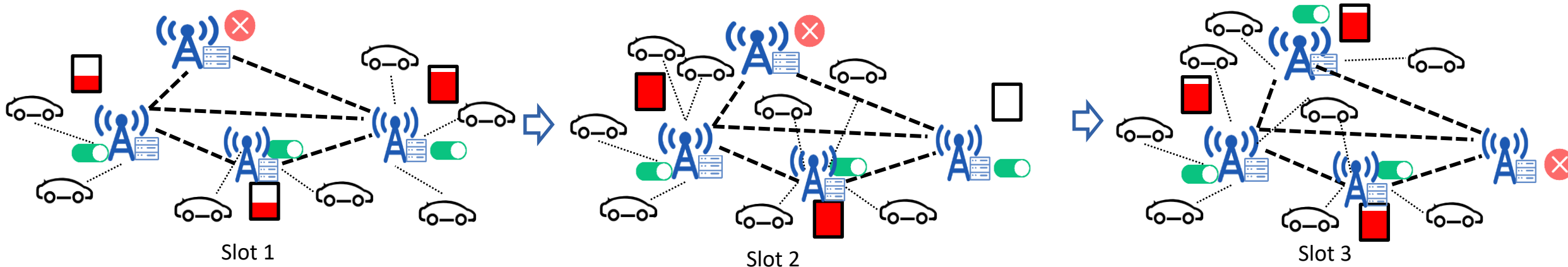


■ Facts observed from the system

- Vehicle distribution **reflects on workload** of edge servers
- Vehicles show **spatio-temporal** relation

■ Basic idea

- **Learn** from the spatio-temporal correlated workload information of edge servers to guide reservation



■ Basic tool

- Deep reinforcement learning



■ RL model

- States: $\mathbf{s}_t = [u_{1,t} + q_{1,t}, \dots, u_{E,t} + q_{E,t}]$

- Actions: $\mathbf{a}_t = [x_{1,t}, \dots, x_{E,t}]$

- Reward: $r_t = \sum_{i=1}^E (-\alpha x_{i,t} + \beta u_{i,t} - \gamma q_{i,t})$

■ Challenge in such an RL problem

- The system has **large state space** (e.g., workload of 1910 edge servers) and **action space** (up to 2^{1910})
- Traditional exploration process is hard to obtain sufficient "good" experience for a DRL agent [1]

[1] Xu, Zhiyuan, et al. "Experience-driven networking: A deep reinforcement learning based approach." in Proc. *IEEE INFOCOM*, 2018.

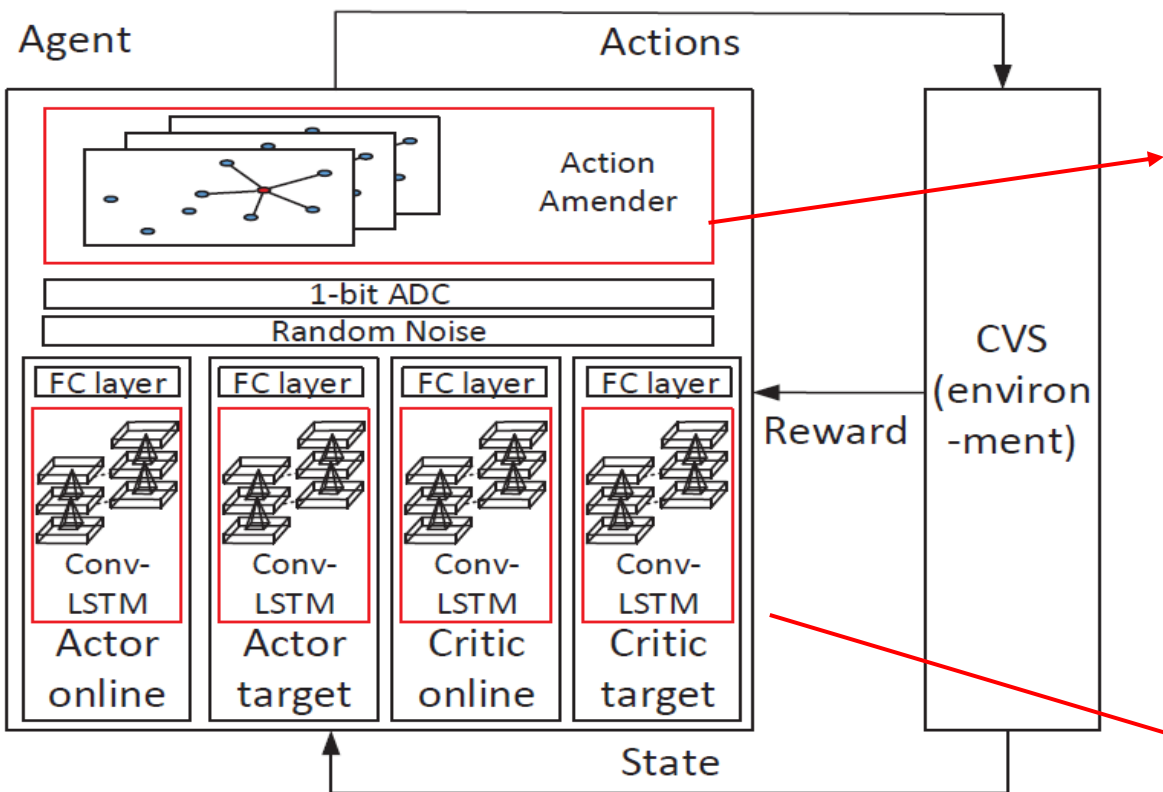


- **Framework: DDPG [1]**
 - Large action space → **policy gradient** → continuous action space
 - Large state space → **deep neural network** → replace large Q tables
 - Speedup training → **Actor Critic**
- **Remaining issues**
 - Fully-connected layers in DNN adopted by DDPG **do not encode spatial and temporal features**
 - output of DNN is not an accurate prediction of future states
 - **Random exploration** from the huge action space is unlikely to gather enough high-reward experiences
 - low system utility during exploration

[1] Lillicrap, Timothy P., et al. "Continuous control with deep reinforcement learning." in Proc. ICLR, 2015.



Two designs to improve DDPG



Framework of DDPG based DRL

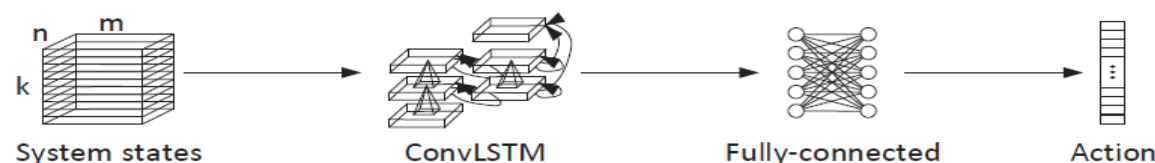
No vehicles in nearby edge servers

No vehicles recently

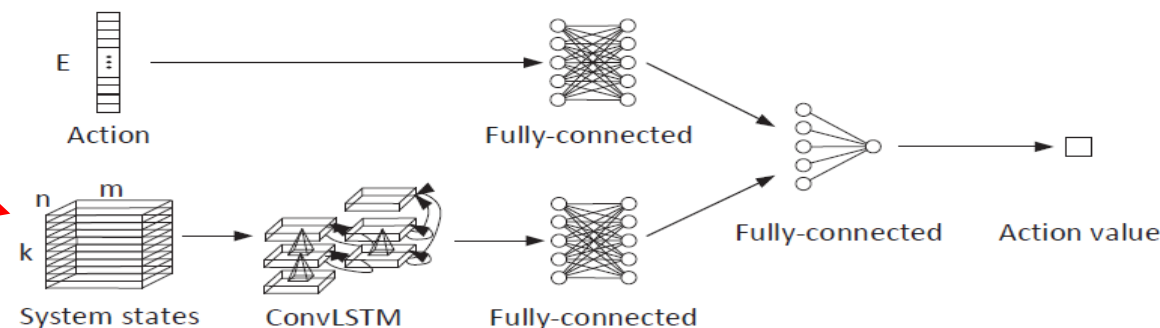
$$X_a: x_{i,t} = x_{i,t}^{(\lambda)} H \left[\sum_{j=0}^l (u_{i,t-j} + q_{i,t-j}) + \sum_{\delta \in \mathcal{E}_g} (u_{\delta,t} + q_{\delta,t}) \right]$$

Network-selected action

Technique 1: **action amender** to modify “bad” actions



(a) The actor network



(b) The critic network

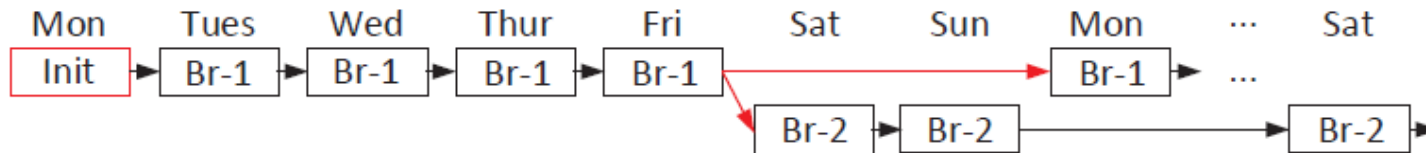
Technique 2: **ConvLSTM** to capture spatio-temporal information

[1] Xingjian, S. H. I., et al. "Convolutional LSTM network: A machine learning approach for precipitation nowcasting." in Proc. **NeurIPS**. 2015.



Two key designs

- Training initializer
 - Apply the greedy algorithm to initialize experience pools
 - Avoid “bad” experiences to deteriorate training
- Forking model branches
 - To adapt to different traffic pattern



Algorithm 1 DR-Train

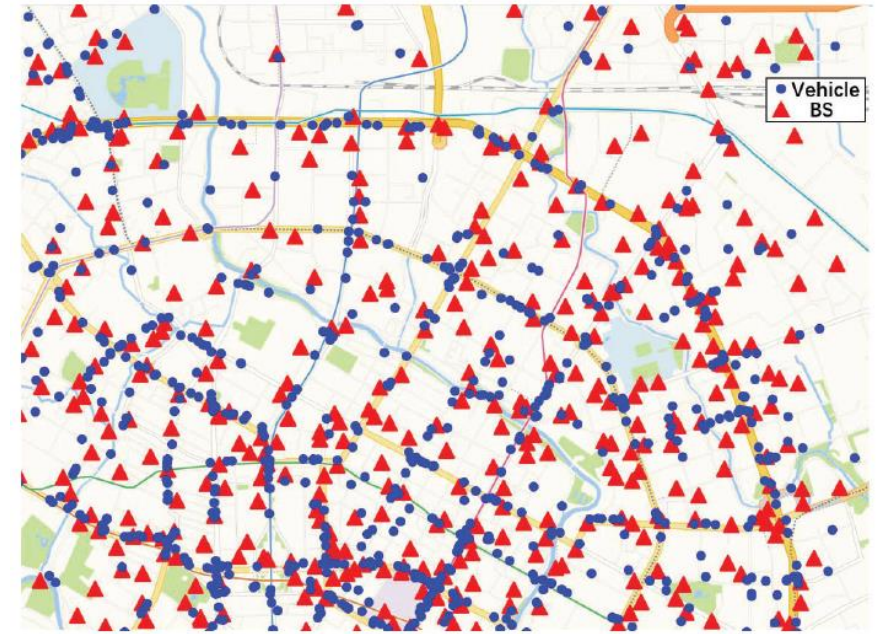
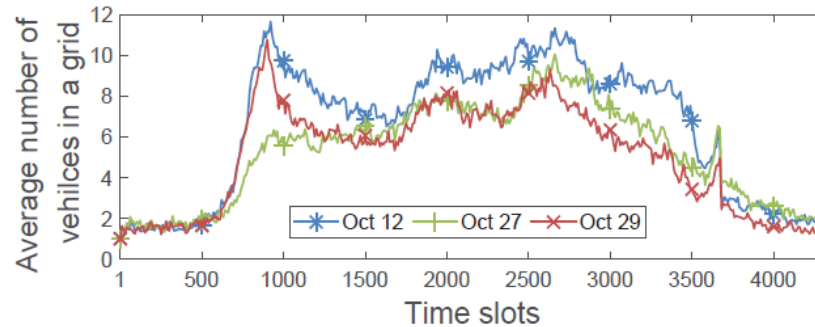
```

1: Randomly initialize critic online network  $Q_w(\cdot)$  and actor
   online network  $\mu_w(\cdot)$  for weekdays with parameters  $\theta_w^Q$ 
   and  $\theta_w^\mu$ , respectively;
2: Initialize target networks  $Q'_w(\cdot)$  and  $\mu'_w(\cdot)$  with parameters
    $\theta_w^{Q'} \leftarrow \theta_w^Q$  and  $\theta_w^{\mu'} \leftarrow \theta_w^\mu$ , respectively;
3: Initialize a random process  $\mathcal{N}$  and experience pools  $R_w$ 
   and  $R_h$  for weekdays and weekends, respectively;
4: Receive initial observation state  $s_1$ ;
5: for  $z$  in  $\mathcal{Z}$  do
6:   if  $z$  is weekday then
7:     for  $t = 1$  to  $T$  do
8:       if  $z$  is the first weekday then
9:         Select  $a_t$  according to the greedy algorithm;
       .....
21:   else
22:     if  $z$  is the first weekend then
23:       Initiate the critic networks  $Q_h(\cdot)$  and  $Q'_h(\cdot)$  and
       .....
    
```



■ Data sets

- Didi express in Chendu in 2018.10.8-11.30
- Base station positions in Chendu



■ Metrics

- training loss of both the critic network and actor network
- system utility
- average resource utilization of reserved MEC servers
- the probability of successful connections among all CVs

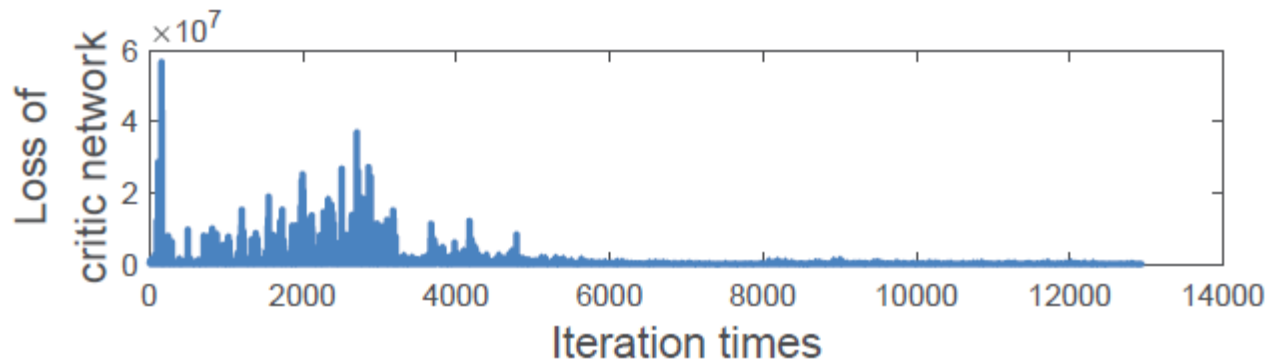
■ Benchmarks

- Variants: 1) DC (remove ConvLSTM), 2) DA (remove action amender), 3) DDPG
- Server placement with demand information: 1) UC [1], 2) HAF [2], 3) GSP [3]
- Others: 1) optimal, 2) random

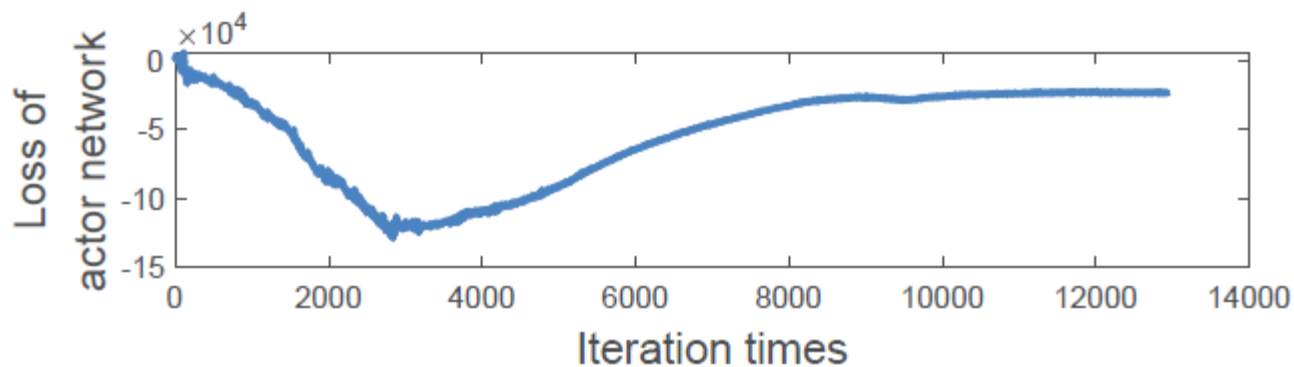
- [1] H. Yin, X. Zhang, H. H. Liu, Y. Luo, C. Tian, S. Zhao, and F. Li, "Edge provisioning with flexible server placement," *IEEE Trans. Parallel Distrib. Syst.*, vol. 28, no. 4, pp. 1031–1045, 2016.
- [2] M. Jia, J. Cao, and W. Liang, "Optimal cloudlet placement and user to cloudlet allocation in wireless metropolitan area networks," *IEEE Trans. Cloud Computing*, vol. 5, no. 4, pp. 725–737, 2015.
- [3] T. He, H. Khamfroush, S. Wang, T. La Porta, and S. Stein, "It's hard to share: Joint service placement and request scheduling in edge clouds with sharable and non-sharable resources," in *Proc. Int. Conf. Distributed Computing Systems (ICDCS)*. IEEE, 2018, pp. 365–375.



Experimental results: DR-Train



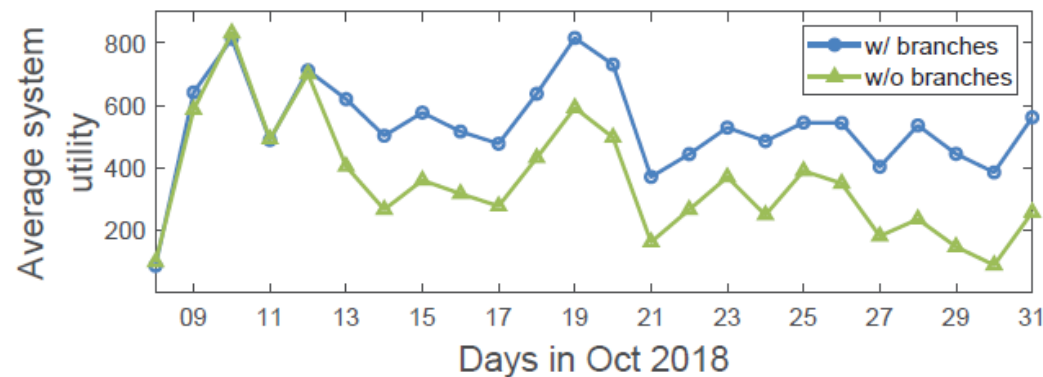
(a) Convergence of the critic network



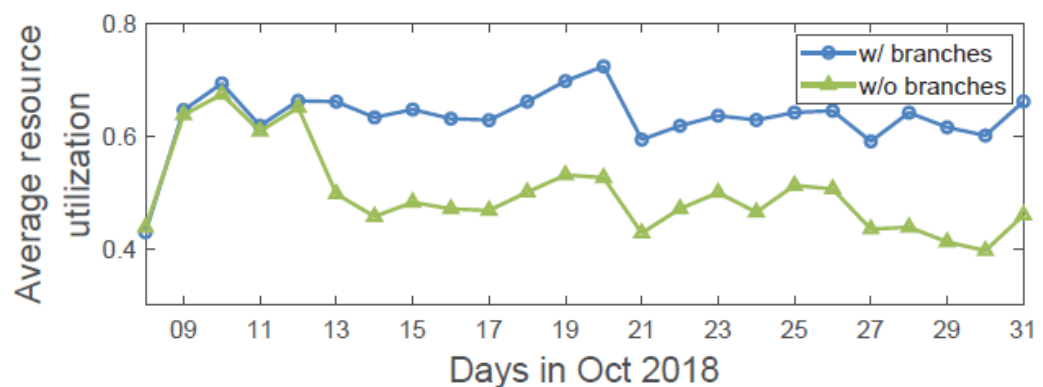
(b) Convergence of the actor network

Fig. 7. The training loss with experience-pool initialization.

Quickly converge to reduce inefficient tries



(a) Average system utility



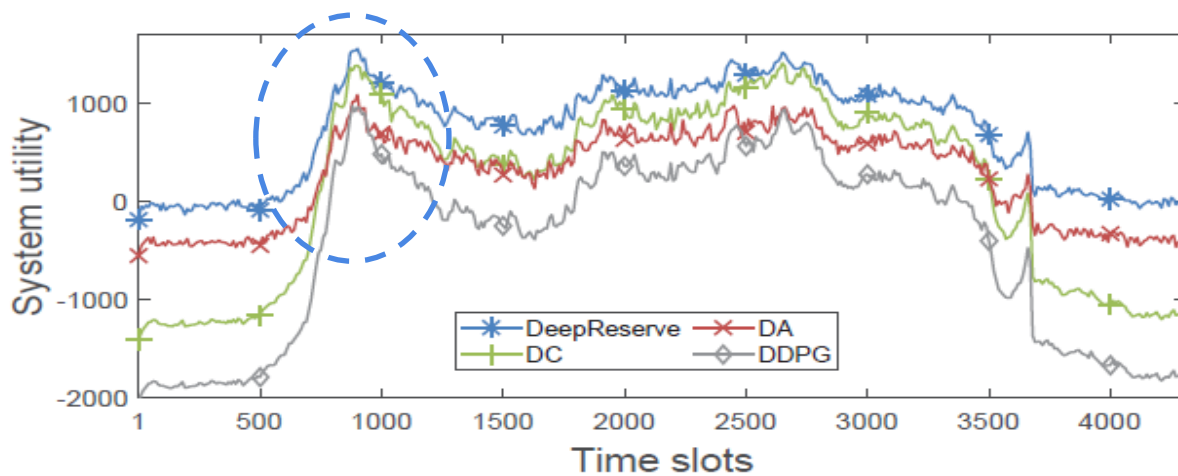
(b) Average resource utilization

Fig. 8. The effectiveness of dividing model branches.

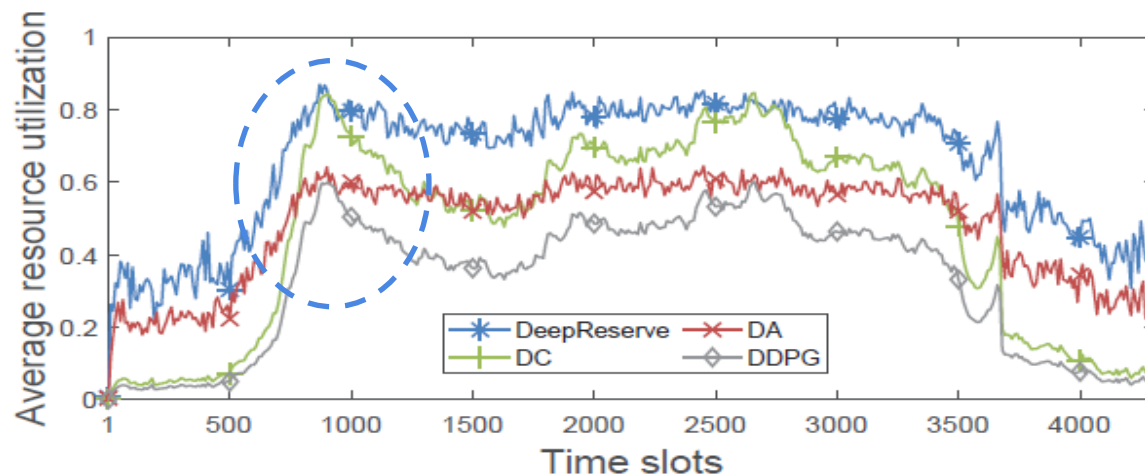
Improve performance with forking branches



Experimental results: variants of DeepReserve



(a) System utility



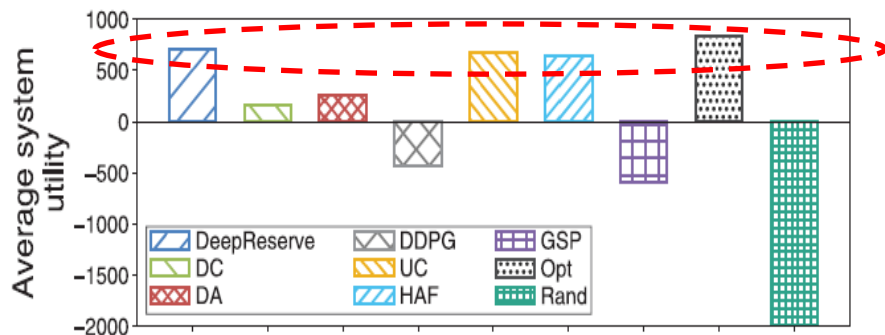
(b) Average resource utilization

Observation: action amender and ConvLSTM contribute differently

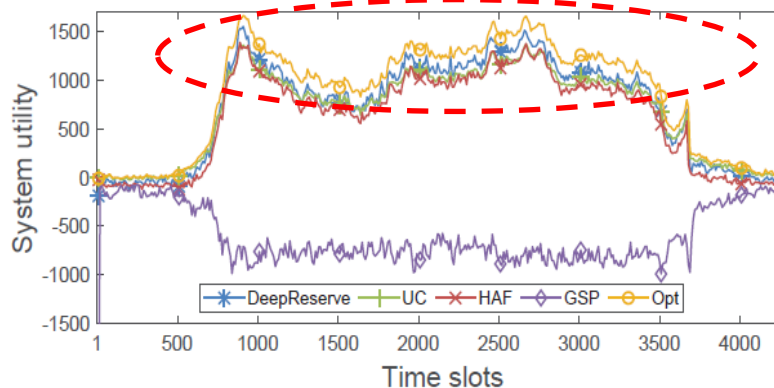
Fig. 11. The performance compared with different variants of DeepReserve.



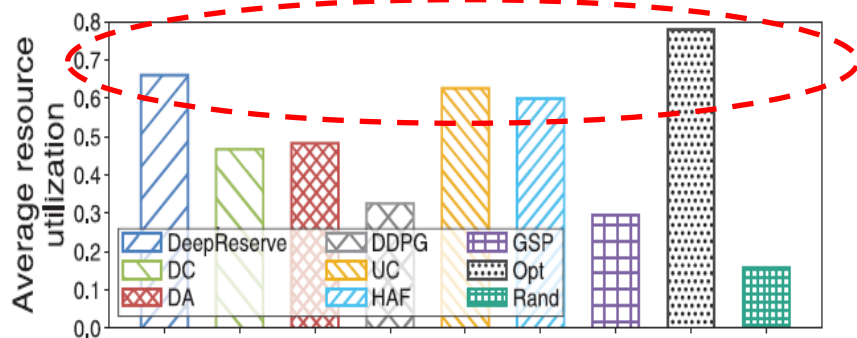
Experimental results: compared to benchmark approaches



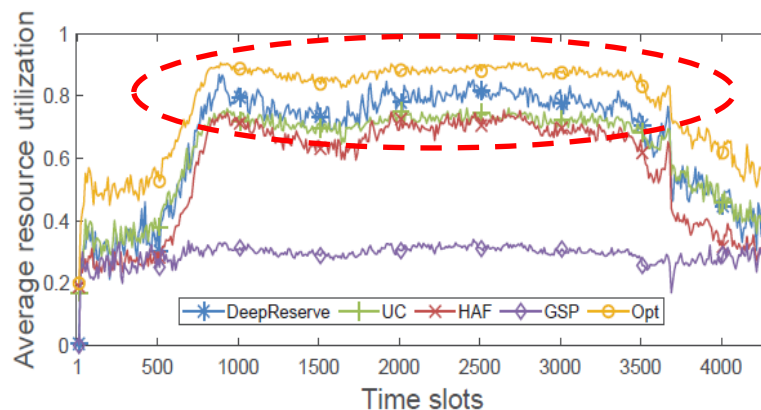
(a) Average system utility



(a) System utility



(b) Average resource utilization



(b) Average resource utilization

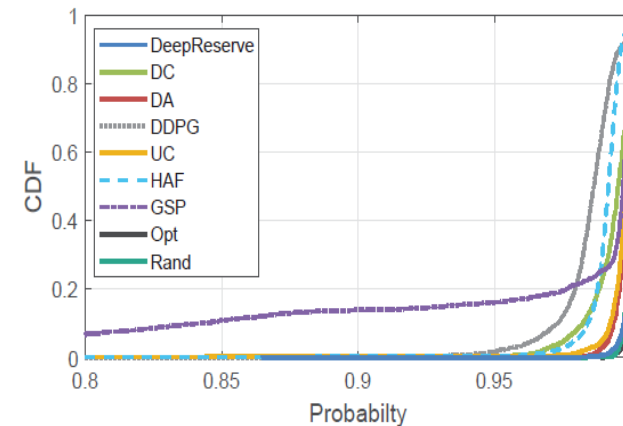


Fig. 13. The CDF plot of the probability of successful connection.

Fig. 10. The performance compared with benchmark approaches. Fig. 12. The performance compared with state-of-the-art approaches.

Near performance with optimal and the approaches with full demand information



■ Contributions

- The **system model** of edge computing based CV system is built and the edge-server reservation problem is **formulated**, which is proved to be NP-hard.
- A DRL based scheme called **DeepReserve** is developed, which is adapted from DDPG with two improvements, i.e., adopting ConvLSTM and the action amender. DeepReserve can efficiently learn to dynamically reserve edge servers without accurate demand information
- A training method called **DR-Train** is designed. Featured with two techniques, i.e., experience-pool initialization and model branches, DR-Train can stably train models for different vehicle traffic patterns.



JOINT INSTITUTE
交大密西根学院

Thank you !

<http://wanglab.sjtu.edu.cn>