

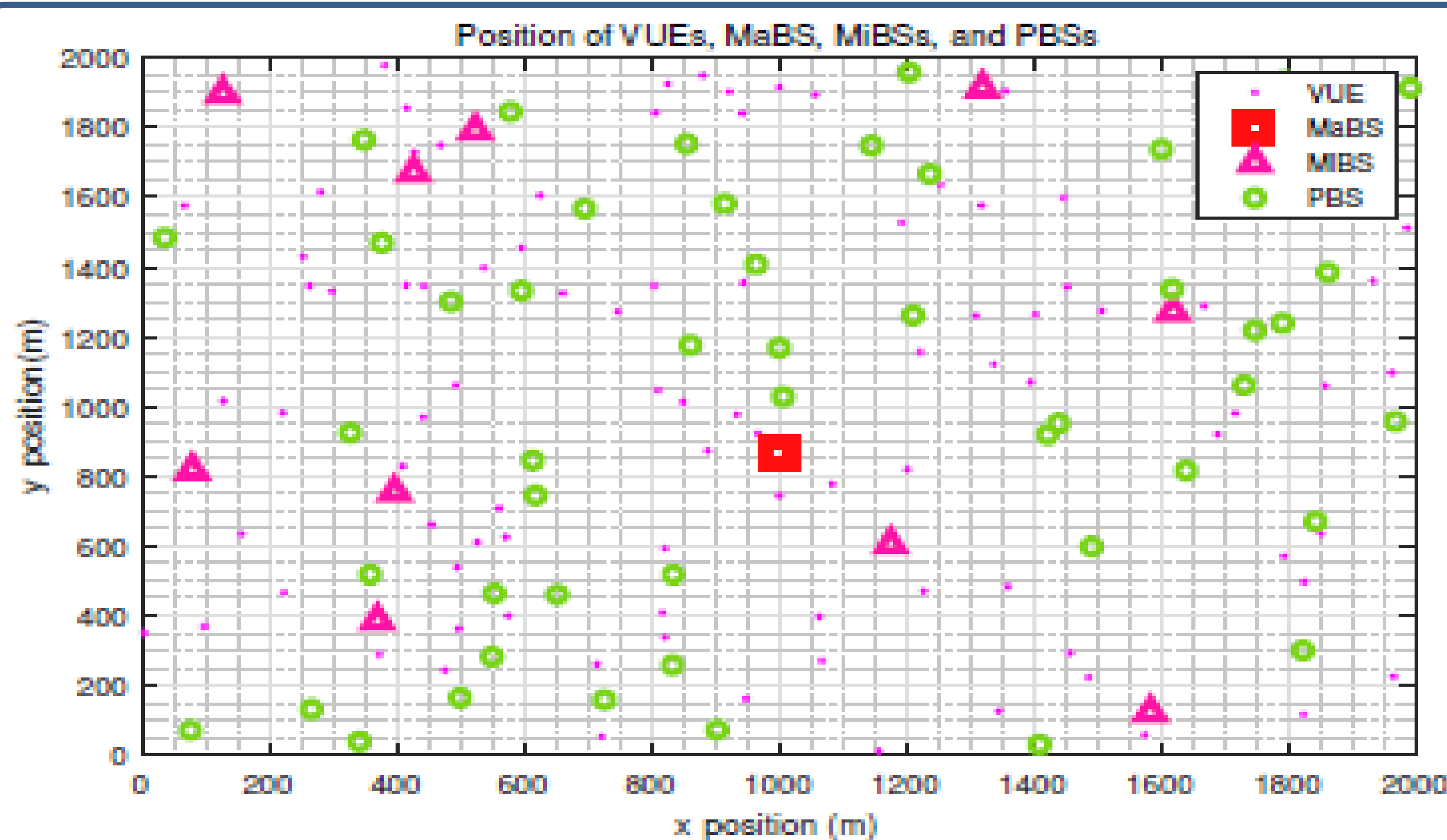
Poster: Multi-Agent Deep Reinforcement Learning for Connected Vehicles

Dohyun Kwon, Soohyun Park, Joongheon Kim (Chung-Ang University, Seoul, Korea)

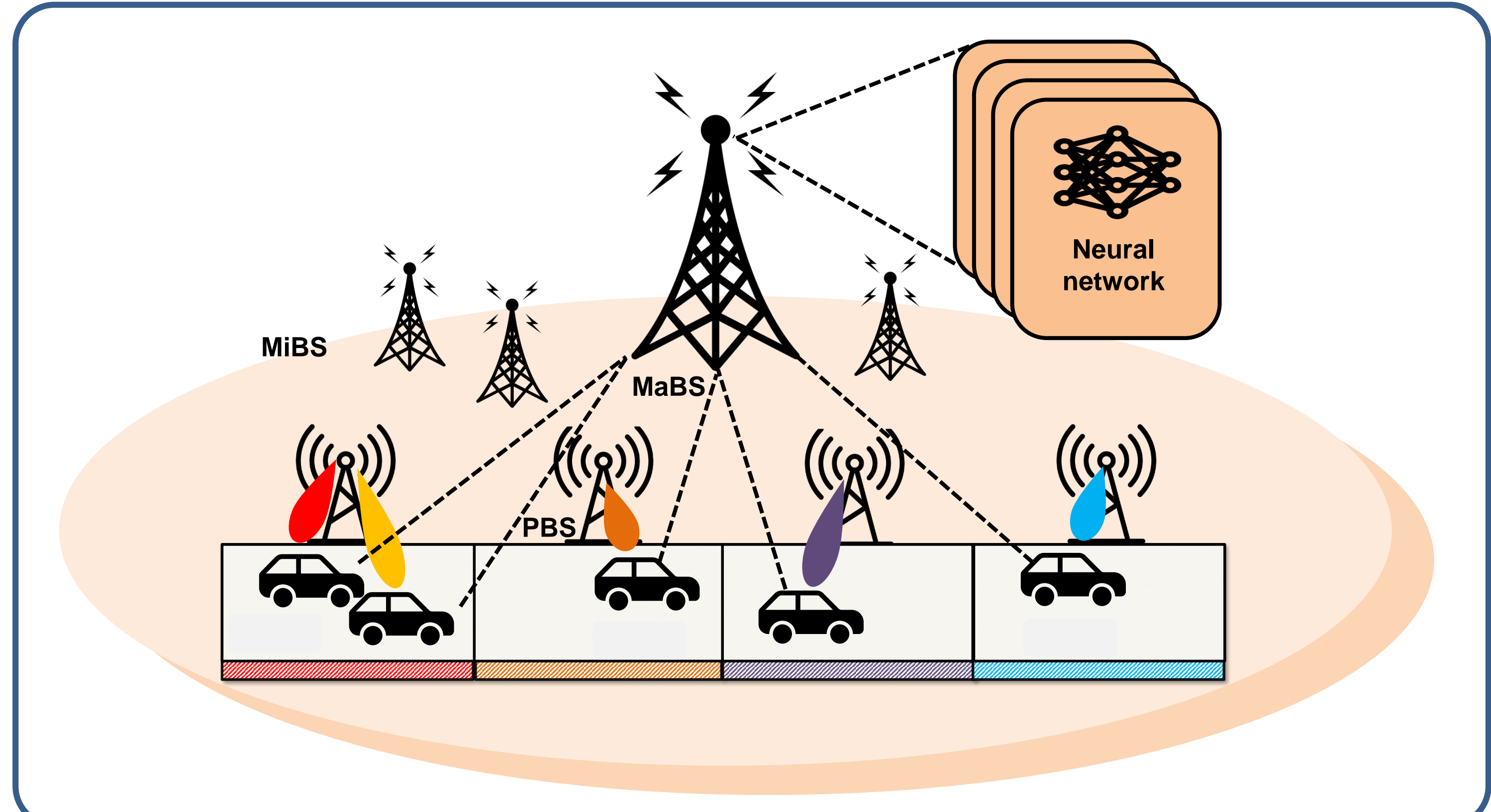
E-mails: kdh1102@cau.ac.kr, shpark.cau@gmail.com, joongheon@cau.ac.kr

Introduction and Backgrounds

Heterogeneous vehicular networks



Connected vehicles in 5G networks



MADDPG for cooperative user association and frequency allocation (UAFA)

State space

- (1) Association information between VUEs and BSs
- (2) Cumulative downlink traffic volume
- (3) Position of VUEs

Action space

- (1) For each VUE, what BS to be associated at next time step
- (2) What channel to select
- (3) If associated BS is PBS, then what direction to set its antenna array to set up the directive link

Reward structure

- (1) High reward when VUEs success to orient its antenna array to a PBS and associate with the PBS
- (2) Small reward for VUEs when they associate with MaBS/MiBS when they can't associate with PBSs
- (3) Penalty for VUEs when they associate with MaBS/MiBS even though they can associate with PBSs

$$Q_{\mu_{\theta}}(\mathbf{x}, a) = r_{t+1} + \gamma \mathbb{E}_{a \sim \mu(\cdot | \mathbf{x}), \mathbf{x} \sim \mathcal{X}} (r_{t+2} + \dots + \gamma^{T-2} r_T). \rightarrow \text{Action-value function}$$

$$\nabla_{\theta_i} \mathcal{J} \approx \frac{1}{S} \sum_j \nabla_{\theta_i} \mu_i(\sigma_i^j) \nabla_{a_i} Q_i^{\mu}(\mathbf{x}^j, a_1^j, \dots, a_N^j |_{a_i = \mu_i(\sigma_i^j)}) \rightarrow \text{Gradient calculation}$$

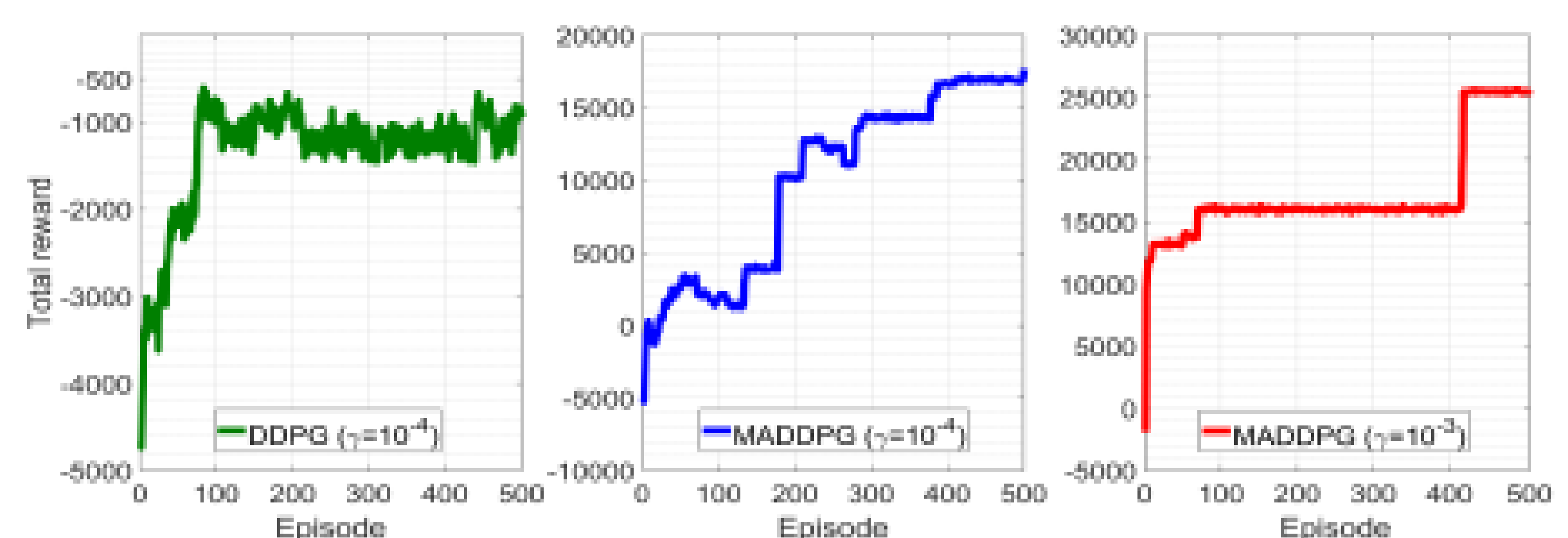
Algorithm description and performance evaluation

Algorithm

Algorithm 1 MADDPG for UAFA in HetVNet

- 1: for episode = 1 to E do
- 2: Initialize the state of VUEs \mathbf{x} and exploration noise \mathcal{N}_t
- 3: for timestep = 1 to T do
- 4: Each i -th selects action $a_i = \mu_{\theta_i}(o_i) + \mathcal{N}_t$
- 5: Execute actions $\mathbf{a} = (a_1, \dots, a_N)$
- 6: Observe r and \mathbf{x}' and store $(\mathbf{x}, \mathbf{a}, r, \mathbf{x}')$
- 7: for VUE $i = 1$ to N do
- 8: Get S samples $(\mathbf{x}^j, a^j, r^j, \mathbf{x}'^j)$ from D
- 9: Set y^j by Eq. (1)
- 10: Update critic by minimizing Eq. (2)
- 11: Update actor by Eq. (3)
- 12: Update θ_i' of each VUE

Performance evaluation



The MADDPG algorithm based UAFA solution showed about 25 times superior performance than single agent based model, which is DDPG based one.