

Heterogeneous Multi-Agent Deep Reinforcement Learning for Traffic Lights Control

Jeancarlo Arguello Calvo

School of Computer Science and Statistics, University of Dublin, Trinity College
arguellj@tcd.ie

Abstract. The most promising technique used in urban traffic control (UTC) has been Reinforcement Learning (RL) due to its capacity to learn the dynamics of complex problems. The main problem of this approach is the curse of dimensionality that arises from the exponential growth of the state and action spaces because of the number of intersections.

Deep Reinforcement Learning (DRL) enhances hugely the performance of RL. DRL has proved to work very well for UTC in single agent environments. Nonetheless, when the problem scales up to multiple intersections the need for coordination becomes too complex, and as a result, the latest studies take advantage of the similarity of agents to train them as one. However, the usage of homogeneous junctions is not a real world scenario where a city has different layouts of intersections.

This study proposes the usage of Independent Deep Q-Network (IDQN) to train heterogeneous multi-agents to deal with both the curse of dimensionality and the need for collaboration. The former is handled by using Deep Q-Networks. The latter is managed with IDQN, which trains simultaneously and separately each agent allowing us to support heterogeneous agents. Unfortunately, this technique can lead to convergence problems because one agent's learning makes the environment appear non-stationary to other agents, and this problem conflicts with experience replay memory on which DQN relies. We address this issue by conditioning each agent's value function on a fingerprint that disambiguates the age of the data sampled from the replay memory.

1 Introduction

This research addresses heterogeneous multi-agents for urban traffic control (UTC) systems as the first study applying Deep Reinforcement Learning (DRL) in a heterogeneous network layout which is a reflection of the cities in the real world. It proposes the usage of DRL to deal with the curse of dimensionality that suffers RL approaches by using Neural Networks. It also uses Independent Q-Learning (IQL), where each agent learns independently and simultaneously its own policy, treating other agents as part of the environment. However, the environment becomes nonstationary from the point of view of each agent, as it involves the interaction with other agents who are themselves learning at the same time, ruling out any convergence guarantees. The technique used for this

problem is Deep Q-Networks which relies in a component so called experience replay memory in order to stabilize and improve the learning. However, this memory is incompatible with non-stationary environments. In order to combine the experience replay memory and IQL, we used a technique of fingerprinting which stabilizes the memory against the non-stationarity. This fingerprint disambiguates the age of the data sampled from the replay memory. We evaluate the usage of IQL with Deep Q-Networks in a multi-agent setting by using a simulation of an urban traffic control system.

2 Background and Related Work

2.1 Reinforcement Learning

Reinforcement learning (RL) is about learning optimal actions from specific situations where the agent has to discover which actions maximize a reward signal by exploring them in a trial-and-error basis. These actions may affect future situations and subsequent reward signals [14].

RL is formulated in the mathematical representation of a complex decision making process called Markov Decision Process (MDP) [13], which it is defined by:

- A time step t .
- A set of states $s \in S$.
- A set of actions $a \in A$.
- A transition function $T(s_t, a_t, s_{t+1})$, which is the probability that an action a leads to s_{t+1} from s_t , i.e. $P(s_{t+1} | s_t, a_t)$.
- A reward function $R(s_t, a_t, s_{t+1})$.

RL is the optimal control of incompletely-known MDP where the transition and the reward functions are unknown. They are learned during the training. In RL an agent interacts with the environment. On each time step t , the agent perceives an state s_t in state space S from where it selects an action a_t in the action space A by following a policy π . The agent receives a reward r_t when it transitions to the state s_{t+1} according to the environment dynamics, the reward function $R(s_t, a_t, s_{t+1})$ and the transition function $T(s_t, a_t, s_{t+1})$. In order to converge the algorithm faster, a discount factor $\gamma \in [0, 1]$ is applied to the total amount of reward. This is defined with the Bellman equation:

$$V^*(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^*(s') \right] \quad (1)$$

2.2 Q-Learning

As in RL, the transition and the reward functions are unknown because it has no prior knowledge of the model (model-free), RL agents learn Q-values instead

of learning values. A Q-value is a function of state-action pair that returns a real value: $Q: S \times A \rightarrow \mathbb{R}$. By using Q-values, the policy can be represented as:

$$\pi(s) = \arg \max_{a \in A} Q(s, a) \quad (2)$$

where the Q-value can be learned by using Q-learning updates [19]:

$$Q(s, a) = (1 - \alpha)Q(s, a) + \alpha[R(s, a) + \max_{a' \in A} Q(s', a')] \quad (3)$$

where $0 < \alpha \leq 1$ is a learning rate.

[19] demonstrated that Q-learning converges to the optimal policy with probability 1 as long as all actions are repeatedly sampled in all states and the Q-values are represented discretely. Exploration is needed to make sure all action in all states are sampled. As a consequence the dilemma of the trade-off between exploration and exploitation arises. A common technique for exploration is to randomly select actions with small probability, this is known as ϵ -greedy exploration.

2.3 Deep Reinforcement Learning

The curse of dimensionality given by large state and action spaces make unfeasible to learn Q value estimates for each state and action pair independently as in normal tabular Q-Learning. Therefore, DRL models the components of RL with deep neural networks. The following subsections describe the methods used during the proposed research to implement the Deep Reinforcement Learning in Traffic Light problem.

Deep Q Networks [11] proposed Deep Q-Networks (DQN) as a technique to combine Q-Learning with deep neural networks. RL is known to be unstable or even to diverge when a non-linear function approximator such as a neural network is used to represent the Q value. DQN addresses these instabilities by using two insights, experience replay and target network.

DQN parameterizes an approximate value function $Q(s, a; \theta_i)$ using Convolutional Neural Networks, where θ_i are the weights of the network at iteration i . The experience replay stores the agent's experiences $e_t = (s_t, a_t, r_t, s_{t+1})$ at each time step t in a dataset $D_t = e_1, e_t$ pooled over many episodes into a replay memory. Then, mini batches of experience drawn uniformly at random from the dataset $(s, a, r, s) \sim U(D)$ are applied as Q-updates during the training. The Q-learning update at iteration i follows the loss function:

$$L_i(\theta_i) = E_{(s,a,r,s) \sim U(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right] \quad (4)$$

where θ_i are the parameters of the Q-network at iteration i and θ_i^- are the target network parameters. The target network parameters are only updated with the Q-network parameters every C steps and are held fixed between individual updates.

Double DQN The max operator in standard Q-learning and DQN uses the same values both to select and to evaluate an action. It is known that this maximization sometimes produces to learn unrealistically high action values which tends to prefer overestimated values over underestimated values, resulting in overoptimistic value estimations. To prevent this, Double Q-learning decouples the selection and the evaluation [7].

For DQN architectures it is not desired to fully decouple the target as in [7] because the target network provides an intuitive option for the second value function, without having to include extra networks. For that reason, [6] propose to evaluate the greedy policy according to the online network, but using the target network to estimate its value given as result the Double DQN (DDQN) algorithm:

$$Y_t^{DDQN} = R_{t+1} + \gamma Q(s_{t+1}, \arg \max_a Q(s_{t+1}, a; \theta_t); \theta_i^-) \quad (5)$$

Prioritized Experience Replay [12] presented prioritized experience replay which detaches agents from considering transitions with the same frequency that they are experienced. Prioritized replay samples more frequently transitions from which there is a high expected learning progress, as measured by the magnitude of their temporal-difference (TD) error. It samples transitions with probability p_t relative to the last encountered absolute TD error:

$$p_t \propto \left| \left(R_{t+1} + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right) \right|^\omega \quad (6)$$

Where ω is a hyper-parameter that determines the pattern of the distribution. New transitions are pushed into the replay buffer memory with maximum priority, providing a bias towards recent transitions.

Dueling Network Dueling Network is a technique proposed by [18] which computes separately the value $V(s)$ and advantage $A(s, a)$ functions that are represented by a duelling architecture that consists of two streams where each stream represents one of these functions. These two streams are combined by a convolutional layer to produce an estimate of the state-action value $Q(s, a)$. The dueling network automatically produces separate estimates of the state value and advantage functions without supervision. Besides that, it can learn which states are valuable, without having to explore the consequence of each action for each state. Dueling network is defined with the equation:

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + \left(A(s, a; \theta, \alpha) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a; \theta, \alpha) \right) \quad (7)$$

2.4 Multi-Agent Reinforcement Learning

Consider a cooperative multi-agent environment where n agents identified by a $\epsilon \in \mathcal{A} \equiv \{1, \dots, n\}$ participate in a stochastic game G , denoted by a tuple $\langle S, U,$

$P, r, Z, O, n, \gamma > 0$. The environment consists of states $s_t \in S$, where at every time step t , each agent takes an action $u_t^a \in U$, forming a joint action $\mathbf{u}_t \in \mathbf{U} \equiv U^n$. The transition probabilities are defined by $P(s_{t+1} | s_t, \mathbf{u}_t) : S \times U \times S \rightarrow [0, 1]$. A global reward function $r(s_t, \mathbf{u}_t) : S \times U \rightarrow \mathbb{R}$ is shared between all the agents.

Each agent's observations $z \in Z$ follows an observation function $O(s_t, \mathbf{a}) : S \times A \rightarrow Z$. Each agent a conditions its behaviour on its own action-observation history $\tau_a \in T \equiv (Z \times U)^*$, according to its policy $\pi_a(u_a | \tau_a) : T \times U \rightarrow [0, 1]$. After each transition, the action u_t^a and new observation $O(s_t, \mathbf{a})$ are added to τ_a , forming τ_a' . Let the joint quantities over agents be in bold, and the joint quantities over agents other than a with the subscript $-a$, so that, e.g., $\mathbf{u} = [u_t^a, \mathbf{u}_{-a}]$.

Independent DQN Independent DQN (IDQN) is an extension of Independent Q-Learning[16] with DQN for cooperative multi-agent environment, where each agent a observes the partial state s_t^a , selects an individual action u_t^a , and receives a team reward, r_t shared among all agents. [15] combines DQN with independent Q-learning, where each agent a independently and simultaneously learns its own Q-function $Q^a(s, u^a; \theta_i^a)$ [3]. Since our setting is partially observable, IDQN can be implemented by having each agent conditioned on its action-observation history, i.e., $Q_a(\tau_a, u_a)$. In DRL, this can be implemented by given to each agent a DQN on its own observations and actions.

As mentioned before, a key component of DQN is the experience replay memory. Unfortunately, the combination of experience replay with IQL appears to be problematic because the non-stationarity introduced by IQL which provokes that the dynamics that generated the data in the agent's experience replay memory no longer indicate the current dynamics in which the agent is learning.

A fingerprint technique was introduced in [4], which states that the disadvantage of IQL is that it ignores the changing over time of the other agents' policies because it perceives the other agents as part of the environment, which causes non-stationarity on its own Q-function. Hence, the Q-function might be made stationary if it conditioned on the other agents' policies.

Based on the results of the study, a fingerprint must be correlated with the true value of state-action pairs given the other agents' policies. It should gradually change over training time in order to allow the model to generalise across experiences in which the other agents execute policies as they learn.

3 Design

This section outlines the approach used to design the different components of the UTC problem as a DRL problem.

3.1 State Representation

The state is a image-like representation of the current state of the simulator environment (Figure 1a), similar to the concept used in [11]. The state consists

of two matrices of 64x64: (1) a binary matrix P for vehicle positions (Figure 1b), and (2) a matrix S for vehicle speeds (Figure 1c). These matrices have been used in previous works such as [5, 9, 10, 17].

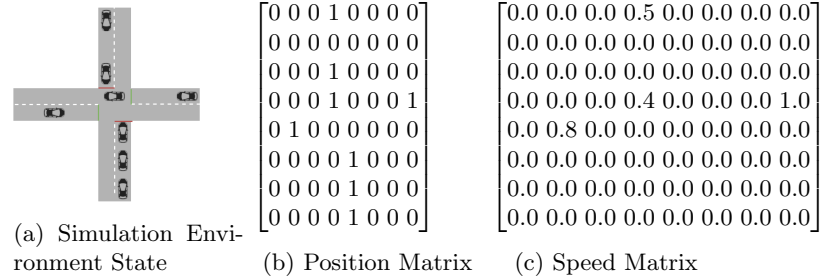


Fig. 1: State Representations

The locations are calculated by mapping the continuous space of the simulated environment into a discretized environment by creating a grid with cells of size C . The matrix P is a binary matrix where one indicates the presence of a vehicle and zero the absence of a vehicle (Figure 1b). The matrix S indicates the speed of the vehicle in the same cell position where a vehicle was calculated to be located in the matrix P (Figure 1c). The speed is represented as a percentage of the maximum allowed speed that it is computed by dividing the current vehicle's speed by the maximum allowed speed.

Additionally, a fingerprint is defined as: Let $-a$ be the other agents of agent a , such that we can include the importance weights of the prioritized experience replay and the TD-Errors vectors sampled from the other agents into the observation function as θ_{-a} and TD_{-a} respectively. Given that, the new observation function is $O'(s) = \{ O(s), \theta_{-a}, TD_{-a} \}$. They are included in a matrix F .

3.2 Action Space

The action space varies depending of the structure of the intersection. For example, in a four road intersection the action space is defined as $A = \{NS, EW, NST, EWT\}$ where NS stands for turning green North-South roads, EW stands for turning green East-West roads, NST stands for turning green North-South right turning, and EWT stands for turning green East-West right turning.

The duration of each phase is 1 time step. However, the agent implicitly determines how long each phase can last ranging from 1 time step up to the final time step of the simulation. The phase can be only changed by the agent when it decides to do it, therefore there is not a limit for how long a phase can take. Additional yellow phase is added during a fixed period of 3 time steps when the previous action is different than the current chosen action. This middle yellow phase reduces the risk of collisions.

3.3 Reward Function

Let $w_{i,t}$ be the i th vehicle's waiting time at time step t , and W_t the total cumulative waiting time for all the vehicles in the observation scope of the road network at time step t as shows in equation 8. The reward function is formulated in equation 9. The intention of the reward function is to reward positive the agent in a range $(0.0, 1.0]$ where the agent's reward loses value proportional to the cumulative waiting time at time step t . Thus, the agent must keep short waiting time in order to receive higher scores, and as consequence, this reward function accomplishes the goal of reducing the driver's waiting time at a junction.

$$W_t = \sum_i w_{i,t} \quad (8)$$

$$r_t = \begin{cases} \frac{1.0}{W_t}, & \text{if cummulative waiting time } W_t \text{ is greater than } 0 \\ 1.0, & \text{otherwise} \end{cases} \quad (9)$$

3.4 Deep Neural Network Architecture

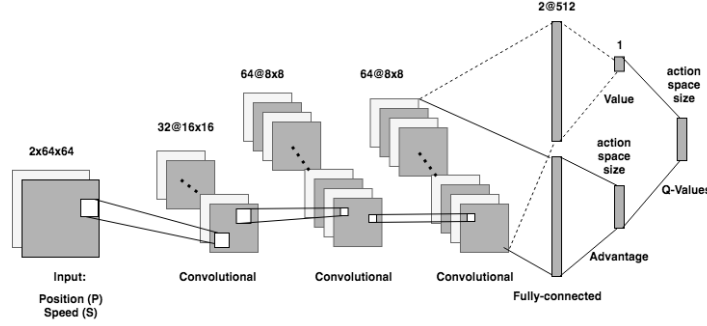


Fig. 2: The architecture of the Deep Neural Network.

The neural network is the same used in [18]. The Figure 2 illustrates the Deep Neural Network. Every agent has its own Neural Network copy in order to allow it to learn its own local policy since every agent is trained independently by using IDQN.

3.5 Framework

We use OpenAI's baseline framework [2] to implement our Dueling Prioritized DDQN. OpenAI Baselines is a set of high-quality implementations of reinforcement learning algorithms in Python which uses the library TensorFlow [1].

4 Evaluation

In this section we present an evaluation of IDQN as solution for multi-agent DRL for heterogeneous agents in traffic light control using the traffic simulator SUMO [8]. We run the test in the network layout show in Figure 3, which consists of three heterogeneous junctions.

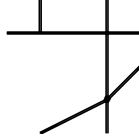


Fig. 3: Network Layout for Multi-Agent Experiments

The table 1 describes the hype-parameters used. Every episode corresponds to 1 hour of simulation time which is 3600 time steps, where every time step is 1 simulated second.

Table 1: Multi-agent evaluation hyper-parameters

Parameter	Value
Episodes	1000
Time steps	3600000
Pre train time steps	2500
Learning rate	0.0004
Exploration ϵ	1.0 \rightarrow 0.01
Time steps from starting ϵ to ending ϵ	360000
Target network update	5000
Prioritization exponent α	0.6
Prioritization importance sampling β	0.4 \rightarrow 1.0
Discount factor γ	0.99
Replay Memory size M	30000
Minibatch size B	32

We select two metrics that have been used in several DRL studies for UTC. Every metric is taken per episode. The metrics are the following:

- **Cumulative Reward.** It is the most common metric for RL experiments due to the performance of the algorithms is directly related with how many rewards the agent can collect over the long time. Bigger numbers means better performance. The reward r_t taken on every time step t is accumulated until the episode is finished.

- **Average Waiting Time.** This metric measures the effectiveness of the reward function into the main goal of the system, which is to reduce the waiting time of every driver in the traffic network. This metric is taken for every i th vehicle that has updated its waiting time w_t^i in the current time step t . Then, the sum of all vehicles' waiting time is the cumulative waiting time W_t at the time step t . When the episode finishes, the sum of all the cumulative waiting time of every time step, from 0 until 3600, is divided by 3600 in order to get the average of the waiting time during an episode as shown in Equation 10. The waiting time is part of the reward function, therefore this metric is correlated with the rewards of every episode. SUMO defines the waiting time for a vehicle as the time (in seconds) spent with a speed below 0.1 m/s since the last time it was faster than 0.1 m/s.

$$AWT = \frac{W_t}{3600} \quad (10)$$

We evaluate the following techniques:

- **Fixed Time / Round Robin (FT).** This is a predefined configuration for all the traffic lights where the duration and order of the phases is fixed and pre configured.
- **IDQN without experience replay (ERM Disabled).** It is a IDQN technique but the experience replay is disabled, therefore the agent cannot store experiences.
- **IDQN with prioritized experience replay (PEMR).** It is a normal IDQN technique which uses the experience replay as recommended for DQN.
- **IDQN with prioritized experience replay and fingerprint (PEMR + FP).** This is the proposed IDQN technique with the fingerprint to disambiguate the age of experience replays.
- **IDQN with standard experience replay (ERM).** This is IDQN with the standard experience replay, not with prioritized experience replay.

Every technique is tested in the following scenarios:

- **Low traffic load.** A scenario where the number of vehicles is reduced. It represents a normal flow where is not peak time. However, it is a good amount of cars to produce traffic jams if the traffic control is not good.
- **High traffic load.** A scenario where the number of vehicles is overwhelming. It represents the traffic load in peak times. This scenario certainly will produced long queues and traffic jams even with good control. The idea is to verify if it is possible and how much can be reduced in such a overcrowded scenario.

4.1 Low Traffic Load Experiments

The Figure 4 presents the set of experiments that we carried out in a low traffic load setup.

Figure 4a shows the cumulative reward obtained per episode for the 4 techniques. FT gets a constant reward range that varies extremely low around 50. This is because this techniques does not adjust its performance, and the variation are a product of the discrepancy of the traffic generation from episode to episode. From this results, we can notice that disabling the experience replay (EMR Disabled) gets very bad rewards, even worse than the baseline, and it never learns. Therefore, disabling the experience replay is not useful for our UTC problem. Finally, we compare the results of the two more promising techniques, i.e. PERM and PERM + FP, out of the four presented. The results shows the both techniques have identical performance getting 400% bigger rewards than FT. They also learn very well in spite of the non-stationarity, increasing the rewards from around 25-100 in the exploration phase up to around 200-250 in the exploitation phase. They reach their optimal point at 300 episodes, where they get to keep a stable rewarding. Based on this similarity on the results of PEMR and PEMR + FP, we deduce that the fingerprint is not helping to get bigger rewards. We analyse the other metrics, and based on that, we conclude why FP is not enhancing the performance.

In terms of waiting time, as expected FT, similarly that in rewards, gets a constant behaviour that varies proportionally to the number of cars generated from each episode. In the case of the EMR Disabled, it shows the direct result of getting very bad rewards which produces extremely high waiting times (results cut in 2000 ms as maximum in order to show the other results in a visible comparable scale). Then, again PEMR and PEMR + FP get alike performance, with very low waiting times which are approximately 300% less than FT, with some episodes with peaks higher than this baseline. From these results we can also conclude that the fingerprinting is not helping to enhance the system.

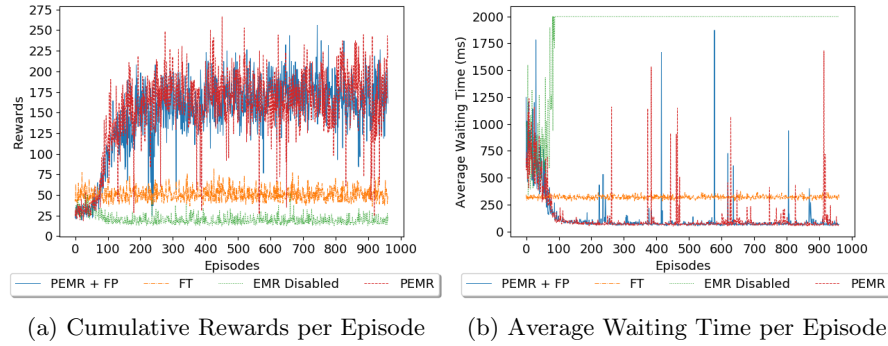
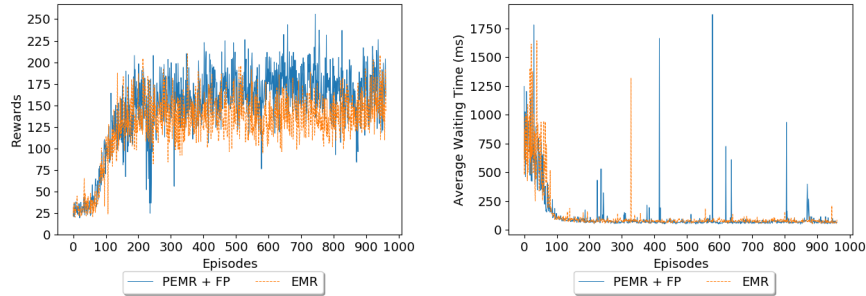


Fig. 4: Low Traffic IDQN Experiments

Now, we move to analyse why the fingerprint is not making better the traffic light control. Our hypothesis holds two possible reasons:

- The components we chose for the fingerprint are not good, because they do not correlate enough with the true value of state-action pairs given the other agents’ policies and/or it does not vary smoothly over training, which does not allow the model to generalise across experiences in which the other agents execute policies as they learn. Likely another component of the DQN training should be used as fingerprint.
- The prioritized experience replay is good enough to deal with the non-stationarity because of main components such as the importance weights it uses to know what records are more valuable samples. This weighting could be enough to disambiguate the age of the data.

Of course, we should run more tests to prove if these hypothesis are true. Because of the long time these experiments take to run, we only tried to verify the second point by executing a test with a standard experience replay. The figure 5 shows this additional experiment which compares the EMR against PEMR + FP.



(a) Cumulative Rewards per Episode (b) Average Waiting Time per Episode

Fig. 5: Low Traffic IDQN Experiments with Standard ERM

As shown in Figure 5a, both techniques, EMR and PEMR + FP, learn well, but PEMR + FP obtains higher rewards outperforming EMR in around 45% in the highest point of PEMR + PR when it gets rewards between 150 and 250 while EMR is getting rewards ranging from 100 to 170. In the case of waiting time illustrated in 5b, they get a similar performance, but EMR is more stable than PEMR + FP by not producing peaks. Regardless of the waiting time, we conclude the Prioritized Experience Replay helps to deal with the non-stationarity by getting bigger rewards, which indicates that PEMR has better performance than EMR. We can minimize the importance of the waiting time results because they are a direct result of the reward function, and PEMR is getting better results from the reward function which is what the agent is going to look for. If by getting bigger rewards the waiting time is not decreasing proportionally, it is a problem with the reward function that do not correlate

that strong with changes in the waiting time, and not a performance problem of the technique. These results prove our second hypothesis that PEMR can be good enough to deal with the non-stationarity. However, we cannot discard the enhancement that fingerprinting can add to PEMR by selecting more appropriate fingerprint components, and therefore, verifying if the first hypothesis is also true. Although, due to time constraint we are not able to demonstrate the first hypothesis.

4.2 High Traffic Load

The Figure 6 presents the set of experiments that we carried out in a high traffic load setup.

The Figure 6a illustrates the results for the cumulative rewards per episode. As can be seen, FT makes a good job in this scenario along with PEMR, both outperforms EMR Disabled and PEMR + FP in around 25%. It seems that the extra layer of the fingerprint reduces the performance of the technique under intense traffic loads. Moreover, none DRL technique gets to learn over the episodes, the performance is constant across all the training. This can indicate that the agents needs more training for such a high load. If we compare the rewarding pattern in Figure 4a, we notice that all the DRL techniques are below FT before episode 100, but after it all PEMR and PEMR + FP beat FT, and EMR Disabled becomes worse. This behaviour might happen after more training episodes in high load. This bad performance in learning can be also produced by a short exploration time, the e-greedy parameter can be increased to allow a more extend exploration of the states.

In terms of the average waiting time, the results are similar to the low traffic load on which techniques perform better. EMR Disabled gets very poor results with a huge magnitude of difference in comparison to all the other techniques. Then, FT gets constant average waiting times as expected reaching around 100% higher times. PEMR and PEMR + FP are the technique which reach lowest waiting times with a similar pattern of results, however PEMR + FP outperform slightly PERM. Unlike the reward results, there is not learning across the training.

5 Conclusion

This paper presented IDQN as the only recent work which addresses heterogeneous multi-agent DRL in UTC systems. We implemented a IDQN system for 3 heterogeneous agent, where agent possesses a Dueling Prioritized DDQN. We evaluated IDQN's performance with different configurations. We executed tests in two traffic conditions, i.e low and high traffic loads. Our experiments showed the IDQN is a suitable techniques for heterogeneous multi-agent UTC environments which can deal with the non-stationarity of the environment. The best results were obtained in the low traffic load. The high traffic load seems to need longer training time. From the successful experiments, we proved the IDQN outperforms the baseline and the experience replay is mandatory to learn

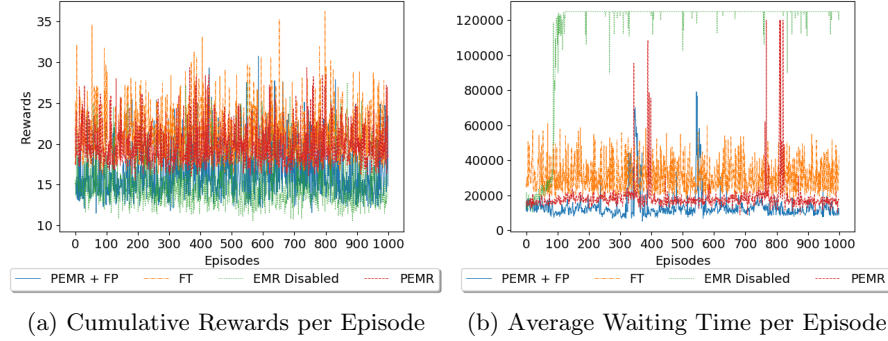


Fig. 6: High Traffic IDQN Experiments

efficiently. Nonetheless, the fingerprint we chose do not add any performance enhancement over the normal prioritized experience replay. Both gets identical performance. We also demonstrated that prioritized experience replay outperforms the standard experience replay in around 45% for the rewards.

Finally, the fingerprint needs further investigation to find out better components as candidates for fingerprint. Also, the way of how we integrate the fingerprint as an additional matrix can be studied to verify if it is the best approach to integrate it into a system which uses a image-like representation as states.

References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D.G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., Zheng, X., Brain, G.: TensorFlow: A System for Large-Scale Machine Learning. In: 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI '16). pp. 265–284 (2016), <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/abadi>
2. Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., Schulman, J., Sidor, S., Wu, Y.: OpenAI Baselines. [\url{https://github.com/openai/baselines}](https://github.com/openai/baselines) (2017)
3. Foerster, J.N., Assael, Y.M., de Freitas, N., Whiteson, S.: Learning to Communicate with Deep Multi-Agent Reinforcement Learning. CoRR abs/1605.0 (2016), <http://arxiv.org/abs/1605.06676>
4. Foerster, J.N., Nardelli, N., Farquhar, G., Torr, P.H.S., Kohli, P., Whiteson, S.: Stabilising Experience Replay for Deep Multi-Agent Reinforcement Learning. CoRR abs/1702.0 (2017), <http://arxiv.org/abs/1702.08887>
5. Gao, J., Shen, Y., Liu, J., Ito, M., Shiratori, N.: Adaptive Traffic Signal Control: Deep Reinforcement Learning Algorithm with Experience Replay and Tar-

- get Network. arXiv pp. 1–10 (2017), <https://arxiv.org/pdf/1705.02755.pdf>{\% }0Ahttp://arxiv.org/abs/1705.02755
6. van Hasselt, H., Guez, A., Silver, D.: Deep Reinforcement Learning with Double Q-Learning. In: Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence. pp. 2094–2100. AAAI’16, AAAI Press (2016), <http://dl.acm.org/citation.cfm?id=3016100.3016191>
 7. Hasselt, H.V., Group, A.C., Wiskunde, C.: Double Q-learning. Nips pp. 1–9 (2010)
 8. Krajzewicz, D., Erdmann, J., Behrisch, M., Bieker, L.: Recent Development and Applications of {SUMO - Simulation of Urban MObility}. International Journal On Advances in Systems and Measurements 5(3&4), 128–138 (dec 2012), <http://elib.dlr.de/80483/>
 9. Liang, X., Du, X., Wang, G., Han, Z.: Deep Reinforcement Learning for Traffic Light Control in Vehicular Networks. Ieee Transactions on Vehicular Technology 1(Xx), 1–11 (2018), <https://arxiv.org/pdf/1803.11115.pdf>
 10. Liu, M., Deng, J., Xu, M., Zhang, X., Wang, W.: Cooperative Deep Reinforcement Learning for Traffic Signal Control. 23rd ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD), Halifax 2017 (2017)
 11. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D.: Human-level control through deep reinforcement learning. Nature 518(7540), 529–533 (2015)
 12. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D., van Hasselt, H., Guez, A., Silver, D., Sorokin, I., Seleznev, A., Pavlov, M., Fedorov, A., Ignateva, A., Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., Hassabis, D., Schaul, T., Quan, J., Antonoglou, I., Silver, D., Rusu, A.A., Gomez Colmenarejo, S., Gulcehre, C., Desjardins, G., Kirkpatrick, J., Pascanu, R., Mnih, V., Kavukcuoglu, K., Hadsell, R., Radford, A., Metz, L., Chintala, S., Ólafsdóttir, H.F., Barry, C., Saleem, A.B., Hassabis, D., Spiers, H.J., Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., Riedmiller, M., Le, Q.V., Ranzato, M., Monga, R., Devin, M., Chen, K., Corrado, G.S., Dean, J., Ng, A.Y., Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D., Ólafsdóttir, H.F., Barry, C., Saleem, A.B., Hassabis, D., Spiers, H.J., Schaul, T., Quan, J., Antonoglou, I., Silver, D., Rusu, A.A., Gomez Colmenarejo, S., Gulcehre, C., Desjardins, G., Kirkpatrick, J., Pascanu, R., Mnih, V., Kavukcuoglu, K., Hadsell, R., van Hasselt, H., Guez, A., Silver, D., Gregor, K., Danihelka, I., Graves, A., Jimenez Rezende, D., Wierstra, D., Kalchbrenner, N., Danihelka, I., Graves, A., Lange, S., Gabel, T., Riedmiller, M., Kalchbrenner, N., Danihelka, I., Graves, A., Gregor, K., Danihelka, I., Graves, A., Jimenez Rezende, D., Wierstra, D.: Prioritized Experience Replay. International Conference in Machine Learning 4(7540), 14 (2015), <http://arxiv.org/abs/1507.01526>{\% }5Cnhttp://dx.doi.org/10.1007/978-3-642-27645-3{_}2{\% }5Cnhttp://arxiv.org/abs/1112.6209{\% }5Cnhttp://arxiv.org/abs/1509.06461{\% }5Cnhttp:

- [//www.arxiv.org/pdf/1509.06461.pdf](http://www.arxiv.org/pdf/1509.06461.pdf)
<http://arxiv.org/abs/1511.06295>
13. Puterman, M.L.: Markov Decision Processes: Discrete Stochastic Dynamic Programming. John Wiley & Sons, Inc., New York, NY, USA, 1st edn. (1994)
 14. Richard S. Sutton and Andrew G. Barto: Reinforcement Learning, Second Edition An Introduction. MIT Press, second edi edn. (2018)
 15. Tampuu, A., Matiisen, T., Kodelja, D., Kuzovkin, I., Korjus, K., Aru, J., Aru, J., Vicente, R.: Multiagent Cooperation and Competition with Deep Reinforcement Learning. arXiv pp. 1–12 (2015), <http://arxiv.org/abs/1511.08779>
 16. Tan, M.: Multi-Agent Reinforcement Learning: Independent vs. Cooperative Agents. In: Machine Learning Proceedings 1993, pp. 330–337. Morgan Kaufmann Publishers Inc. (1993), <http://linkinghub.elsevier.com/retrieve/pii/B9781558603073500496>
 17. Van Der Pol, E., Oliehoek, F.A.: Coordinated Deep Reinforcement Learners for Traffic Light Control. NIPS'16 Workshop on Learning, Inference and Control of Multi-Agent Systems (Nips) (2016), <http://www.fransoliehoek.net/docs/VanDerPol16LICMAS.pdf>
http://elisevanderpol.nl/papers/vanderpol_oлиеhoek_nipsmlc2016.pdf
 18. Wang, Z., de Freitas, N., Lanctot, M.: Dueling Network Architectures for Deep Reinforcement Learning. arXiv (9), 1–16 (2016), <http://arxiv.org/abs/1511.06581>
 19. Watkins, C.J.C.H., Dayan, P.: Q-learning. Machine Learning 8(3-4), 279–292 (1992), <http://link.springer.com/10.1007/BF00992698>