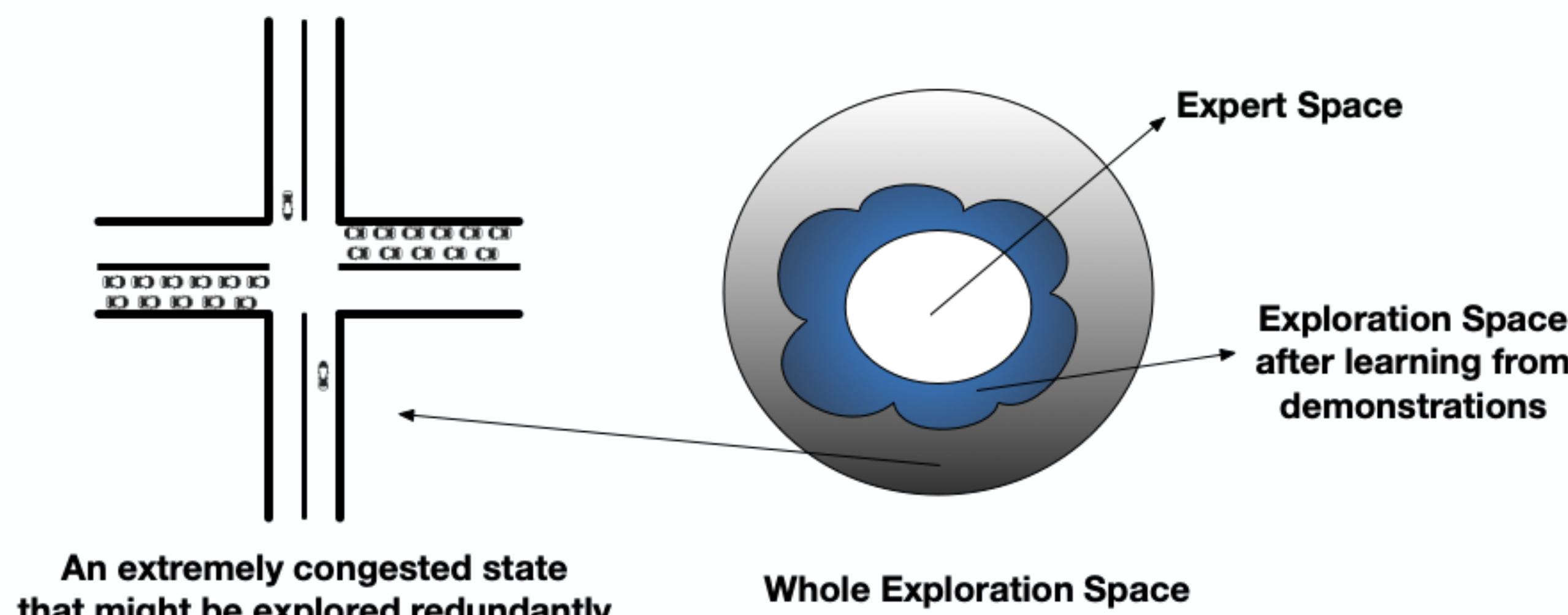


## 1. Introduction

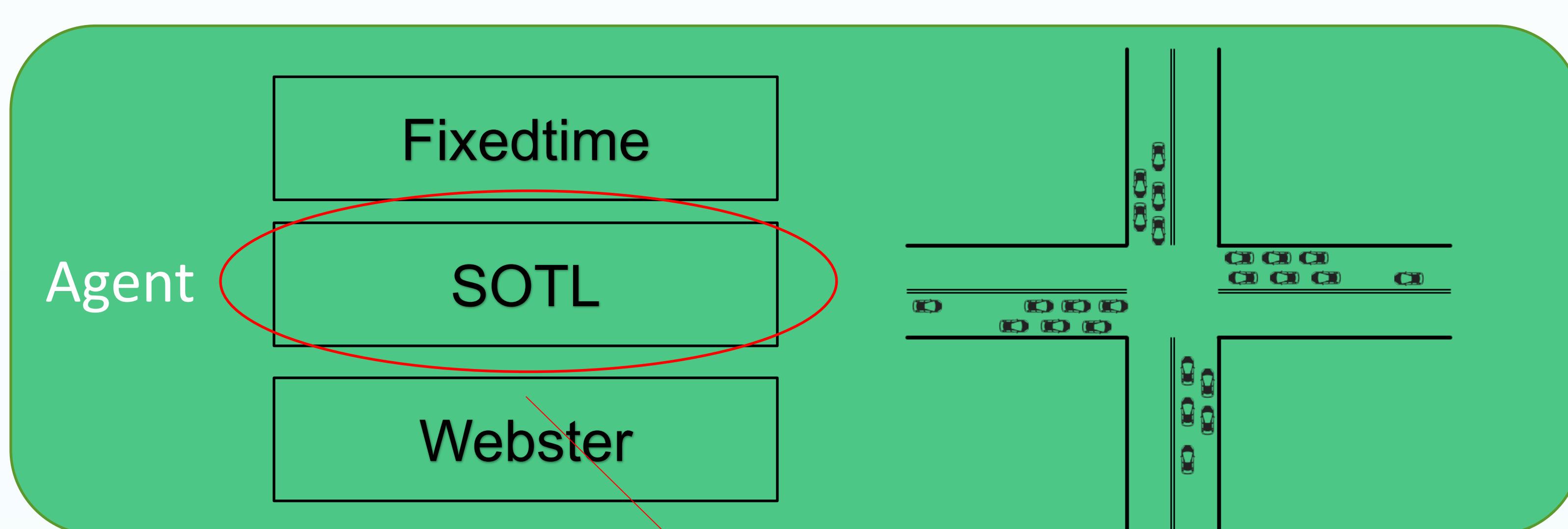
### Large exploration space



## 2. Why learning from demonstrations

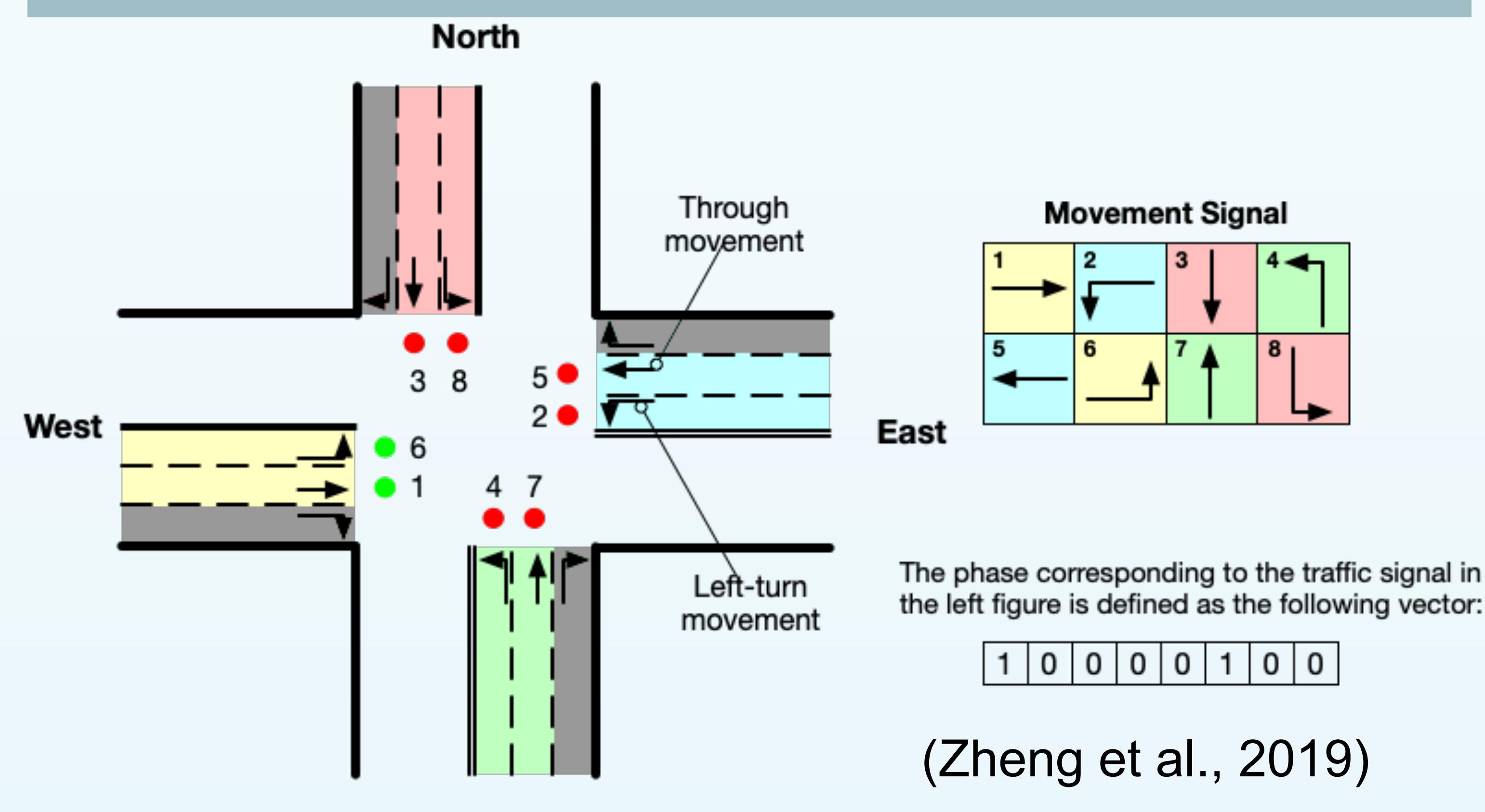
### Experts

### Imitation



Self-Organizing Traffic Lights(SOTL) is used as the expert in our method. It controls the signal by calculating the number of waiting vehicles. Once the number exceeds a threshold, the signal becomes green for those vehicles to pass.

## 3. Setting



- State:** #vehicles on each traffic movement, current phase
- Action:** the phase for the next time interval
- Reward:** average queue length

We are implementing our model in Hangzhou China!

Try to find more related researches? Just scan QR code on the right or visit <https://traffic-signal-control.github.io>

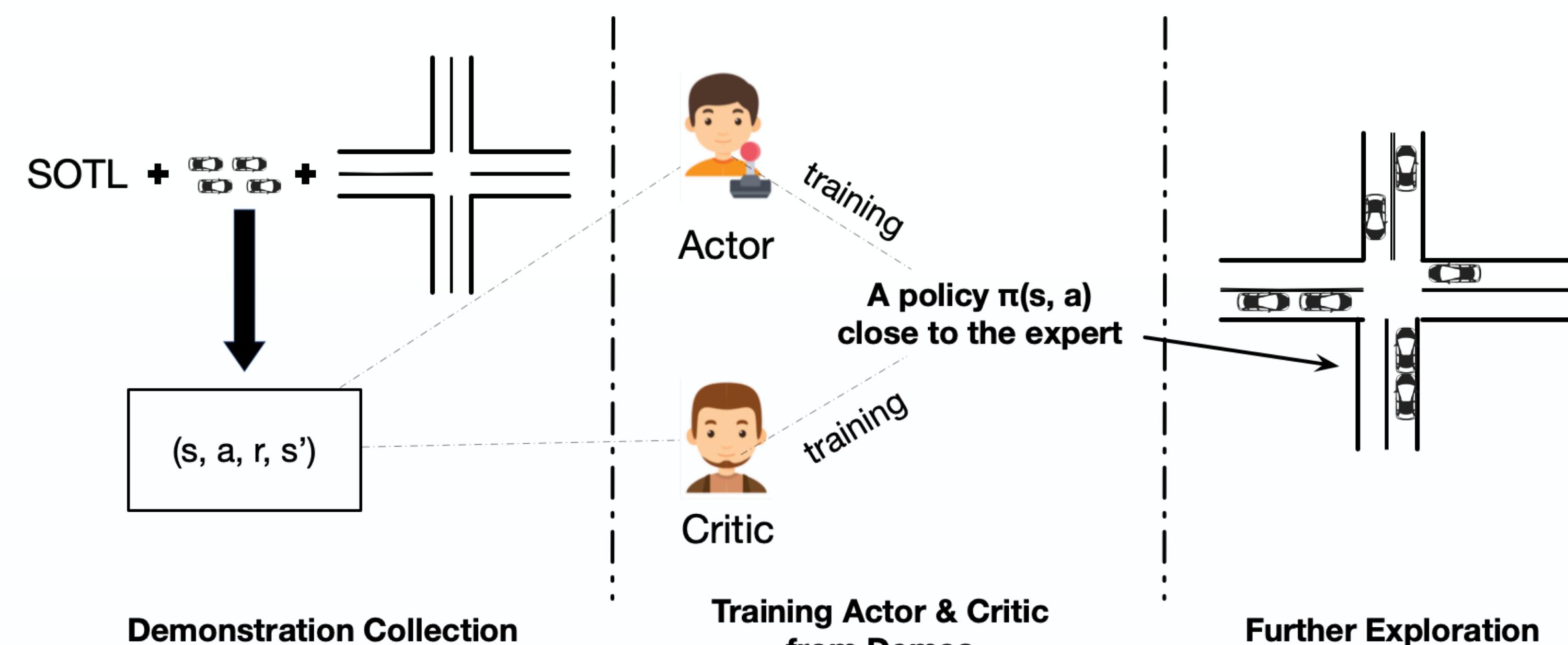


## References

- Wei et al., IntelliLight: A Reinforcement Learning Approach for Intelligent Traffic Light Control
- Wei et al., A Survey on Traffic Signal Control Methods
- Zheng et al., Diagnosing Reinforcement Learning for Traffic Signal Control
- Wei et al., PressLight: Learning Max Pressure Control to Coordinate Traffic Signals in Arterial Network
- Zheng et al., Learning Phase Competition for Traffic Signal Control

- Wei et al. , CoLight: Learning Network-level Cooperation for Traffic Signal Control
- Xiong et al., Learning Traffic Signal Control from Demonstrations
- Zhang et al., CityFlow: A Multi-Agent Reinforcement Learning Environment for Large Scale City Traffic Scenario
- Hester et al., Deep Q-learning from demonstrations

## 4. Methodology



Figures of actor and critic are from <https://www.freecodecamp.org/news/an-intro-to-advantage-actor-critic-methods-lets-play-sonic-the-hedgehog-86d6240171d/>

- To make the policy network differentiable

$$a_{soft} = \text{softmax}((g_i + \pi)/\tau)$$

- With the action of demos  $a_D$  as the ground truth, the loss function can be derived:

$$L_{demo}(\theta_\pi) = \text{Cross-Entropy}(a_{soft}, a_D)$$

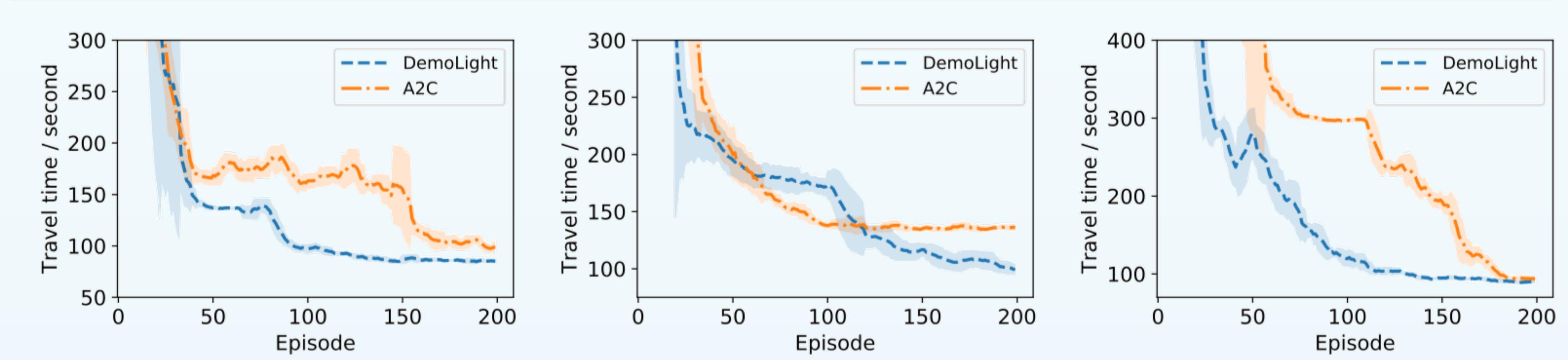
- Following Google's DQfD (Hester et al., 2018), the loss of the value network is computed as:

$$L_{demo}(\theta_Q) = L_{1-TD}(\theta_Q) + L_{n-TD}(\theta_Q) + L_{margin}(\theta_Q) + L_{reg}(\theta_Q)$$

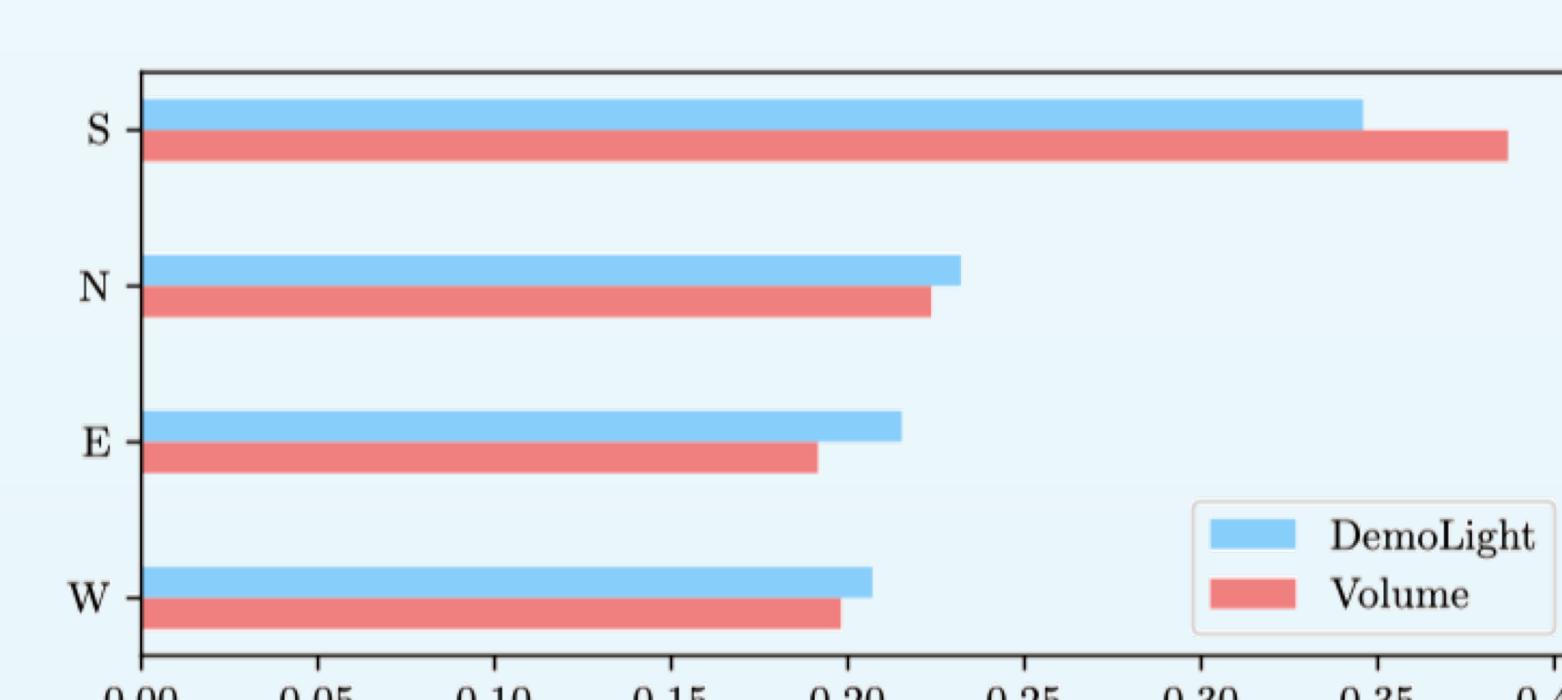
## 5. Experiment

Table 1: Overall performance. Travel time is reported in the unit of second.

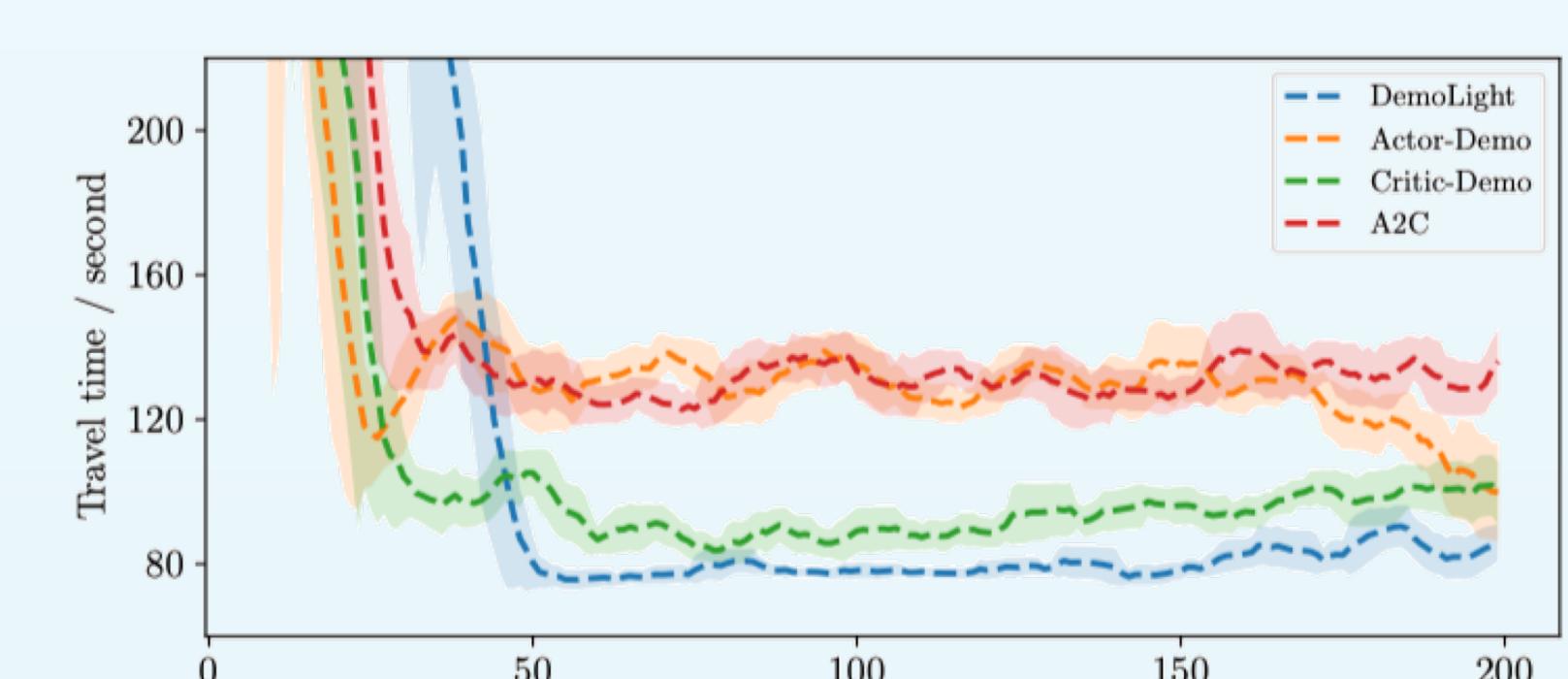
Model	City A			City B			City C				
	1	2	3	1	2	3	4	5	1	2	3
SOTL [2]	102.76	179.41	248.73	248.12	153.43	165.38	123.73	269.64	89.36	149.49	72.11
LIT [16]	103.31	122.88	154.77	346.30	146.88	139.83	104.84	551.65	93.95	130.11	48.43
A2C [1]	85.48	98.36	136.01	380.64	92.14	135.93	77.56	517.77	76.38	96.74	46.83
DemoLight	<b>76.64</b>	<b>85.16</b>	<b>85.93</b>	<b>116.10</b>	<b>92.02</b>	<b>97.88</b>	<b>76.39</b>	<b>183.70</b>	<b>72.44</b>	<b>91.26</b>	<b>43.04</b>
Improvement	10.34%	13.42%	36.82%	69.50%	0.13%	27.99%	64.52%	1.98%	5.30%	5.66%	8.09%



### Convergence Speed



### Learnt Policy



### Ablation Study

## 6. Acknowledgement

The work was supported in part by NSF awards #1652525, and #1618448. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.