

# Asynchronous Distributed Reinforcement Learning for Traffic Signal Control<sup>1</sup>

Shalabh Bhatnagar

CSA Department

Joint Affiliation: RBCCPS and CiSTUP

IISc

December 16, 2016

---

<sup>1</sup>Joint work with Prabuchandran K.J and Hemanth Kumar H.N. 

# Outline of the Presentation

- Introduction and setting
- Basic Framework
- Our algorithm
- Numerical results
- Conclusions and Future work

# Traffic Signal Control (TSC) problem

- The problem - intelligently control traffic lights at junctions to reduce the average delay experienced by the road users
- Possible control scenarios:
  - A Centralized control<sup>23</sup>
  - B Decentralized control<sup>4</sup>
- We control the green time duration of phases by fixing an order for the phase sequence.

---

<sup>2</sup>Prashanth L.A and S.Bhatnagar, *IEEE Transactions on Intelligent Transportation Systems*, 12(2):412-421, 2011

<sup>3</sup>Prashanth L.A and S.Bhatnagar, *IEEE Transactions on Vehicular Technology*, 61(9):3865-3860, 2012

<sup>4</sup>Prabuchandran K.J, Hemanth Kumar A.N, and S.Bhatnagar, *IEEE ITSC*, pp.2529-2534, 2014

## TSC Problem

Find the optimal green signal duration of phases at each intersection to minimize average delay of road users

### Existing approaches

- Fixed time signal control - Not adaptive
- Manual signal control - Locally good actions
- Automated adaptive signal control - Learning approaches

# Available approaches

- Fixed Timing Control : Signal times calculated offline using traffic characteristics at different times of day<sup>5</sup>
- Adaptive Signal Controls :
  - Model Based : Use mathematical models for traffic dynamics such as fluid patterns of traffic<sup>6</sup>
  - Model Free : Reinforcement Learning algorithms<sup>7</sup>

---

<sup>5</sup>M.Papageorgiou, *Proceedings of the IEEE*, 91:2043-2067, 2003

<sup>6</sup>N.H.Gartner, C.J.Messer, and A.K.Rathi, *U.S. Department of Transportation, Transportation Research Board*, 1992

<sup>7</sup>B.Abdulhai, R.Pringle and G.Karakoulas, *Journal of Transportation Engineering*, 129:278-290, 2003

# Available Information

- We assume that some information about traffic congestion across various lanes is available through sensors
- The available information need not be precise but could be 'coarse' in nature
- We assume the phase-switching in an intersection happens in a round robin manner
  - Ensures fairness
  - Ideally a passenger would like to know which signal will turn green next

# Centralized Control Scenario

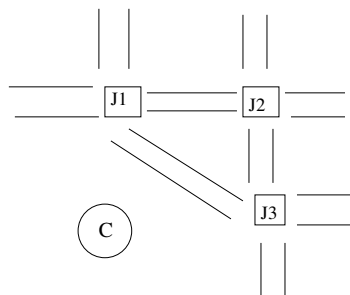
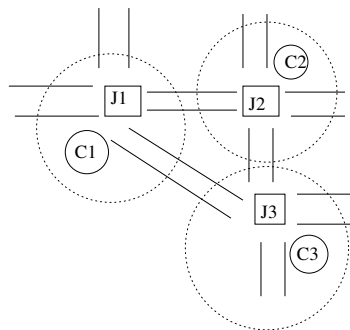


Figure: A Centralized Control Setting – Single Controller

- Controller receives congestion information from all junctions
- Post computation sends feedback on signal configurations to each junction

# Decentralized Control Setting



**Figure:** A Decentralized Control Setting – Multiple Controllers

- Each controller makes decisions using local state information
- Information exchange across junctions factored in local computations



# Our approach

- We formulate TSC problem as a co-operative multi-agent system whose goal is to minimize average delay
- Each junction is modelled as an agent and each agent handles the state space associated with it
- Each agent uses Q-learning with either  $\epsilon$ -greedy<sup>8</sup> or UCB based<sup>9</sup> exploration
- The feedback cost signal to an agent is obtained from itself and its neighbours

---

<sup>8</sup>R.Sutton and A.Barto, *Reinforcement Learning*, Cambridge University Press, 1998

<sup>9</sup>P.Auer, N.Cesa-Bianchi, and P.Fischer, *Machine Learning*, 47(2):235-256, 2002

Multiagent asynchronous learning with each agent (junction) exchanging congestion information with its one-hop neighbors

## Design Issues

- Modeling the problem of optimizing signal duration of phases in a junction with real time traffic - MDP Framework
- Low response time - Coarse information on traffic conditions
- Scalability of the approach - Distributed learning

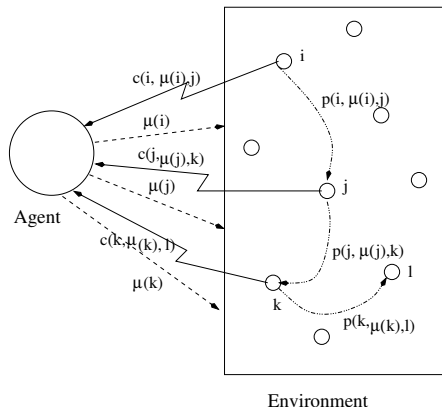


Figure: Agent-Environment Interaction

<sup>10</sup>M.L.Puterman, *Markov Decision Processes*, Wiley, 1994

<sup>11</sup>D.P.Bertsekas and J.N.Tsitsiklis, *Neuro-Dynamic Programming*, Athena Scientific, Belmont, MA, 1996

# The Goal

- Find a policy  $\mu^*$  that gives the optimal cost

$$V^*(i) \equiv V_{\mu^*}(i) \triangleq \min_{\mu} E \left[ \sum_{j=0}^{\infty} \gamma^j c(s_j, \mu(s_j), s_{j+1}) \mid s_0 = i \right], i \in S$$

- $V^*$  satisfies the Bellman equation

$$V^*(i) = \min_{a \in A(i)} \sum_j p(i, a, j) (c(i, a, j) + \gamma V^*(j))$$

# The Q-Value Function

- Let  $Q^*(i, a)$  be defined according to

$$Q^*(i, a) = \sum_j p(i, a, j)(c(i, a, j) + \gamma V^*(j))$$

- Then

$$V^*(i) = \min_{a \in A(i)} Q^*(i, a)$$

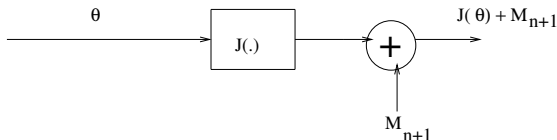
- Moreover,

$$Q^*(i, a) = \sum_j p(i, a, j)(c(i, a, j) + \gamma \min_{b \in A(j)} Q^*(j, b))$$

- The above is also known as the Q-Bellman equation

# Stochastic Approximation<sup>12</sup>

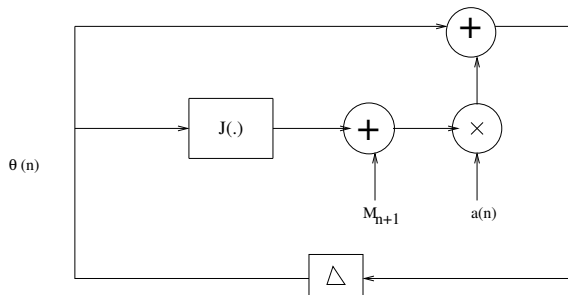
- **Objective:** Solve the equation  $J(\theta) = 0$  when analytical form of  $J$  is not known, however, noisy measurements  $J(\theta(n)) + M_{n+1}$  can be obtained
- Here  $M_n$ ,  $n \geq 0$  constitutes the 'noise', zero-mean and uncorrelated with the past



<sup>12</sup>H. Robbins and S. Monro *Annals of Mathematical Statistics*, 22: 400–407, 1951

# The Robbins-Monro Algorithm

- Algorithm  $\theta(n+1) = \theta(n) + a(n)(J(\theta(n)) + M_{n+1})$



- $a(n) = \frac{1}{n}, \frac{1}{n^\alpha}, \alpha \in (0.5, 1], \frac{\log(n)}{n}, \frac{1}{n \log(n)}$  etc.

# More General Recursions

- Consider the recursion

$$\begin{aligned}\theta(n+1) &= \theta(n) + a(n)h(\theta(n), \mu(n)) \\ &= \theta(n) + a(n)(J(\theta(n)) + M_{n+1}),\end{aligned}\tag{1}$$

where  $J(\theta(n)) = E[h(\theta(n), \mu(n)) \mid \theta(n)]$

- Under some assumptions, algorithm has the same asymptotic behaviour as

$$\dot{\theta}(t) = J(\theta(t))$$

- Important instances:

- Fixed point iteration:  $J(\theta) = f(\theta) - \theta$
- Stochastic gradient scheme:  $J(\theta) = -\nabla f(\theta)$



# An SA Algorithm for Q-Bellman Equation

- Recall the Q-Bellman equation

$$Q^*(i, a) = \sum_j p(i, a, j)(c(i, a, j) + \gamma \min_{b \in A(j)} Q^*(j, b))$$

- Let

(1)  $\theta = Q(i, a)$

(2)  $J(\theta) = \sum_j p(i, a, j)(c(i, a, j) + \gamma \min_{b \in A(j)} Q(j, b)) - Q(i, a)$

- Goal: Find  $\theta^* \triangleq Q^*(i, a)$  such that  $J(\theta^*) = 0$

- An iterative RL algorithm which solves Q-Bellman equation using following update rule:

$$\begin{aligned} Q_{n+1}(i, a) &= Q_n(i, a) + a(n)(c(i, a, j) + \gamma \min_{b \in A(j)} Q_{n+1}(j, b) - Q_n(i, a)) \\ &= Q_n(i, a) + a(n)(J(\theta(n)) + M_{n+1}) \end{aligned}$$

where  $a(n)$  : step size at time 'n'

---

<sup>13</sup>C.J.C.H.Watkins and P.Dayan, *Machine Learning*, 8(3):279-292, 1992

- ODE associated with Q-learning:

$$\dot{Q}(i, a) = \sum_j p(i, a, j)(c(i, a, j) + \gamma \min_{b \in A(j)} Q(j, b)) - Q(i, a)$$

- The above ODE has  $Q^*(i, a)$  as its unique globally asymptotically stable equilibrium
- Then  $Q_n(i, a) \rightarrow Q^*(i, a)$  almost surely as  $n \rightarrow \infty$
- *Problem:* The procedure is computationally expensive and requires huge memory

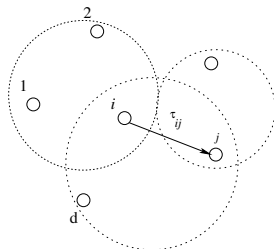


Figure: Asynchronous Updates

- Let  $\theta_i(n)$  be  $n$ th parameter update of  $i$ th processor
- Each processor updates its parameter component on its own 'local' clock

<sup>14</sup>S.Bhatnagar, *Systems and Control Letters*, 60:472-478, 2011

# Stability and Convergence

- The Distributed Stochastic Update: For  $i = 1, \dots, d$ ,

$$\theta_i(n+1) = \theta_i(n) + a(\nu(i, n))I\{i \in Y_n\}$$

$$(J_i(\theta_{1-\tau_{i1}}(n), \dots, \theta_{d-\tau_{id}}(n)) + M_{n+1}(i))$$

- Here  $\nu(i, n) = \sum_{m=0}^n I\{i \in Y_m\}$  and  $Y_n$  is the (random) subset of components updated at time  $n$

- Key requirement:  $\liminf_{n \rightarrow \infty} \frac{\nu(i, n)}{n} > 0$  for all  $i$

- Key result:  $\sup_n \|\theta(n)\| < \infty$  almost surely and  $\{\theta(n)\}$  asymptotically tracks the ODE

$$\dot{\theta}(t) = \Lambda(t)J(\theta(t))$$

## State space of MDP

- $L_j$  - number of incoming lanes to a junction  $j$ .
- State of junction  $j$  at time ' $t$ ' is defined as a vector

$$s^j(t) = (q_1^j(t), q_2^j(t), \dots, q_{L_j}^j(t), p^j(t))$$

$q_i^j(t)$  - queue length in the  $i^{\text{th}}$  lane at time  $t$

$p^j(t)$  - index of the phase that has to be set to green

- State space of the network:  $S \triangleq \times_{j=1}^J S^j$   
 $J$  - total number of junctions and  $S^j$  state space of junction  $j$
- This poses a problem as state space increases rapidly with number of junctions in the network

# Solution: Decentralization

## State space

- We consider each junction as a separate learning agent
- Each agent observes only a portion of the state space
- This consideration improves scalability

## State space of a junction

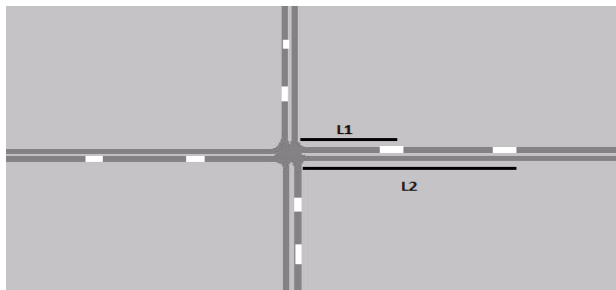
- Consider a junctions with 4 incoming lanes, each of which accommodates 100 vehicles
- Total possible states for the junction -  $4 * 10^8$
- State space for a junction is still considerably large



# Solution: State Discretization

Segregate the queue-length in a lane as  $\{\text{low}=0, \text{medium}=1, \text{high}=2\}$  using sensors

$$q_i^j(t) = \begin{cases} \text{low}, & \text{if } q_i^j(t) < L1 \\ \text{medium}, & \text{if } L1 \leq q_i^j(t) \leq L2 \\ \text{high}, & \text{if } q_i^j(t) > L2 \end{cases}$$



# Action space and Cost function

- **Actions :**

- Action  $a_j$  of agent  $j$  in state  $s_j$  is to decide the green light duration of current phase
- The action set is  $\{low = 10s, medium = 20s, high = 30s\}$  same for all states

- **Cost Function :** Cost  $c_j(t)$  at time  $t$  incurred by agent  $j \in J$  for taking an action  $a_t^j$  at state  $s_t^j$  :

$$c_j(s_t^j, a_t^j, s_{t+1}^j) = \frac{1}{|N_j|} \sum_{k \in N_j} \sum_{i=1}^{L_k} q_i^k(t+1), \quad (2)$$

$N_j$  - set including itself and neighbors of junction  $j$

# Multiagent Reinforcement-Learning (MARL) Algorithm

- Each agent updates its Q-factor using the Q-learning update rule
- Agent at junction  $j$  follows :

$$Q_{n+1}^j(s^j, a^j) = Q_n^j(s^j, a^j) + a(n) (c_j(n) + \gamma \min_{b \in A} Q_n^j(s^{j'}, b) - Q_n^j(s^j, a^j))$$

$s^j, s^{j'}$  - queue-length vectors at junction  $j$  at times  $n$  and  $n + 1$

action  $a^j$  at time  $n$  is set according to one of the exploration mechanisms ( $\epsilon$ -greedy exploration or UCB exploration)

- The cost feedback signal  $c_j(n)$  is obtained from neighbours for taking action  $a^j$

## $\epsilon$ -greedy exploration strategy

With probability  $1 - \epsilon$ , choose action  $a^j = \arg \min_c Q_n(s^j, c)$  and with probability  $\epsilon$ , choose random feasible action  $a^j$  in state  $s^j$

## UCB-based exploration strategy

Compute action  $a^j$  based on upper confidence bound index

$$a^j = \arg \max_{c \in A} \left\{ -Q_n^j(s^j, c) + \sqrt{\frac{\ln R_{s^j}(n)}{R_{s^j,c}(n)}} \right\}, \quad (3)$$

$R_{s^j}(n)$  - Number of visits to state  $s^j$  until time  $n$

$R_{s^j,c}(n)$  - Number of times action  $c$  is chosen in state  $s^j$  until time  $n$

If action  $c$  is not chosen sufficient number of times in state  $s^j$  then second term will dominate and action  $c$  will be explored. With learning, first term dominates and action is selected only based on Q-factors.

---

## Algorithm 1

---

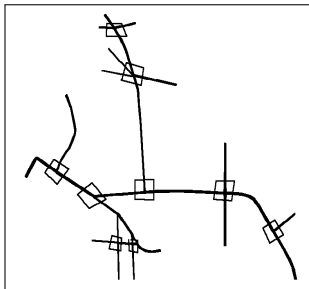
```
for  $time = 1$  to  $\infty$  do
  for  $j = 1$  to  $number\_of\_junctions$  do
    if  $green\_time(j) = 0$  then
      a.  $curr\_phase(j) = (curr\_phase(j)+1) \bmod number\_of\_phases(j)$ 
      b. Set  $green\_time(j)$  using one of the exploration mechanisms
      c. Collect local and neighbor's traffic conditions to compute cost
      d. Update Q-factors for junction  $j$  using Q-learning update rule
    end if
     $green\_time(j) = green\_time(j) - 1$ 
  end for
end for
```

---

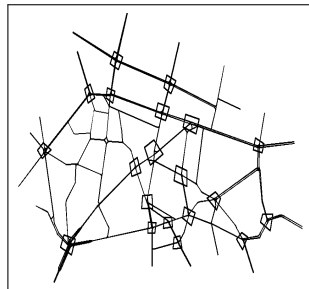
# Experimental Setup

- Simulations are run on VISSIM, a microscopic traffic simulator
- Two different road networks
  - A nine-junction road network (9Jn) - a real road network around the Indian Institute of Science campus in Bangalore, India.
  - A twenty-junction road network (20Jn)- a real road network in Bangalore, India.

# Road Network Layouts



(a) Nine-Junction Network



(b) Twenty-Junction Network

Table: Different road network layouts



# Parameter Settings

- Our action set is  $\{low = 10s, medium = 20s, high = 30s\}$
- Traffic conditions for the road network

Road network	Arrival rate per hour
9 Jn	5500
20 Jn	9650

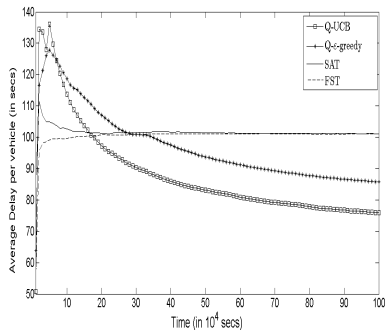
# Algorithms Considered

- We compare performance of our MARL algorithms with Fixed Signal Timing (FST) and Saturation Balancing (SAT)<sup>15</sup> algorithms
  - FST algorithm - Signal timing tuned for best performance
  - SAT algorithm - an adaptive algorithm that controls signal timings based on the vehicle density during last cycle
- Performance metrics - average delay and average stopped delay

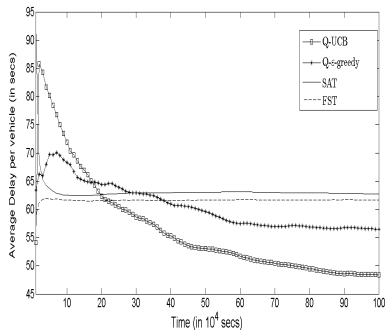
---

<sup>15</sup>S.Richter, D.Aberdeen and J.Yu, *NIPS*, 2007

# Performance Comparison - Average Delay

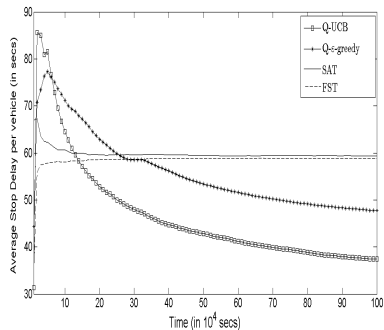


(a) 20 Jn road network

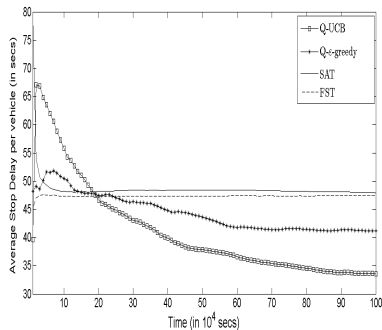


(b) 9 Jn road network

# Performance Comparison - Average Stopped Delay



(a) 20 Jn road network



(b) 9 Jn road network

# Performance Statistics

	Average stopped delay [s]	Average delay [s]	Average number of stops
9JnFST	47.40	61.72	2.14
9JnSAT	47.92	62.74	2.11
9JnQ- $\epsilon$ -greedy	41.13	56.46	1.82
9JnQ-UCB	33.60	48.41	1.63
20JnFST	58.78	101.06	4.03
20JnSAT	59.30	101.10	4.01
20JnQ- $\epsilon$ -greedy	47.61	85.74	3.56
20JnQ-UCB	37.31	75.82	3.04

**Table:** Average performance comparisons for the road network layouts

- The threshold parameters in the policies can be tuned adaptively for optimal performance<sup>16</sup>
- Non-zero communication delays between junctions need be considered
- Function approximation approaches can be explored together with asynchronous distributed techniques
- Signalling for pedestrian movement
- Modeling driver behaviour

---

<sup>16</sup>Prashanth L.A and S.Bhatnagar, *IEEE Transactions on Vehicular Technology*, 61(9):3865-3880, 2012