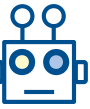


Machine Learning for Networking

# Reinforcement Learning

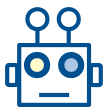
Session 10 – Q-learning II

Boris Bellalta: [boris.bellalta@upf.edu](mailto:boris.bellalta@upf.edu)

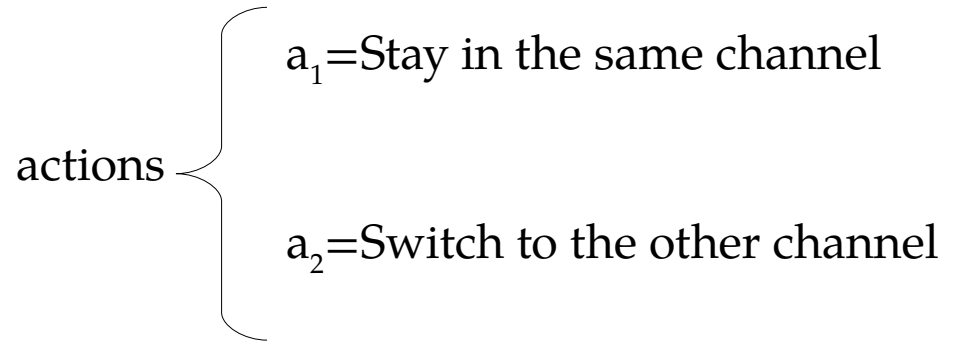
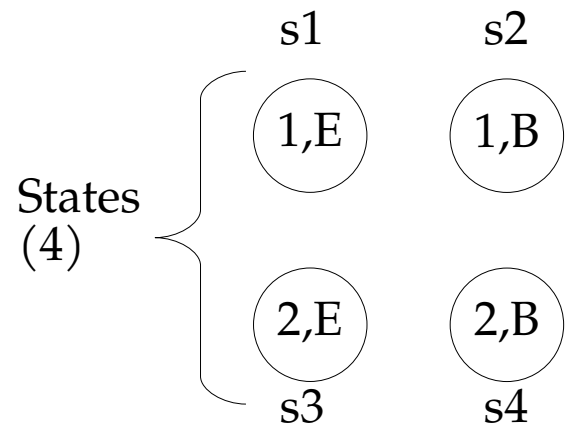
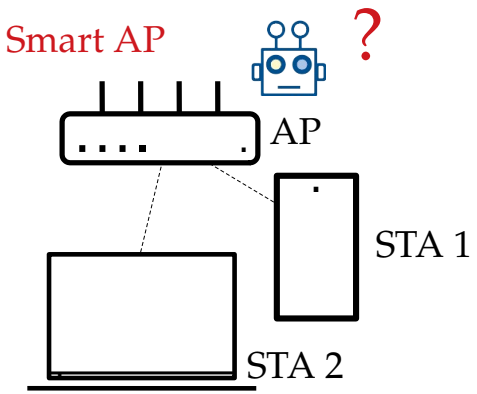


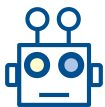
# Contents

- Q-learning: step by step, 2<sup>nd</sup> part
- Q-learning exercises
- Exercise: video congestion control → What is the effect of rewards?



# States & Actions





# Q-learning

- Table to store the pair (state,action), i.e. the Q-table

State	Stay	Switch
1,E (s1)		
1,B (s2)		
2,E (s3)		
2,B (s4)		

- A mechanism to explore the state space

- Epsilon-greedy (for example)

- A way to update the Q-table (from wikipedia) – Bellman's equation

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \left( \underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)$$

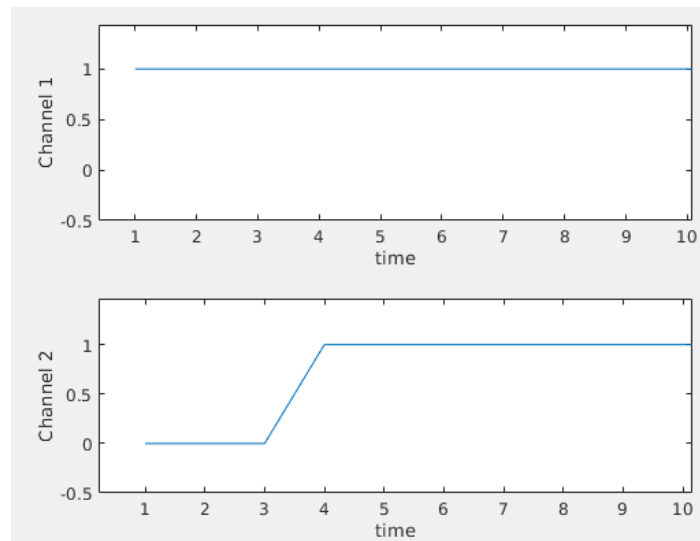
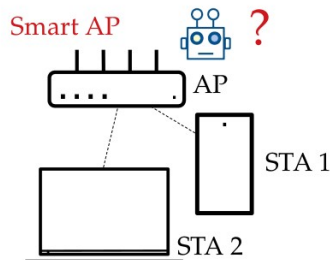
temporal difference

new value (temporal difference target)

# Case I

channels 1 2

$\alpha=0.2$   
 $\gamma=0.9$



I have no info at the start

```
0.0000e+000  0.0000e+000
0.0000e+000  0.0000e+000
0.0000e+000  0.0000e+000
0.0000e+000  0.0000e+000
```

Iteration | Current State

```
1.0000e+000  2.0000e+000
```

Explore? | Action (1 = stay | 2 = switch)

```
0.0000e+000  1.0000e+000
```

Next state | Reward (1 = well done | -1 = :( )

```
2.0000e+000  -1.0000e+000
```

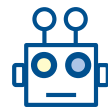
$(1-\alpha) * Q\_table(state(t),action(t)) + \alpha * (reward(t) + \gamma * \max(Q\_table(state(t+1),:))) =$

```
800.0000e-003  0.0000e+000  200.0000e-003  -1.0000e+000  900.0000e-003  0.0000e+000
```

-----Q-table-----

```
0.0000e+000  0.0000e+000
-200.0000e-003  0.0000e+000
0.0000e+000  0.0000e+000
0.0000e+000  0.0000e+000
```

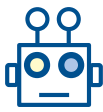
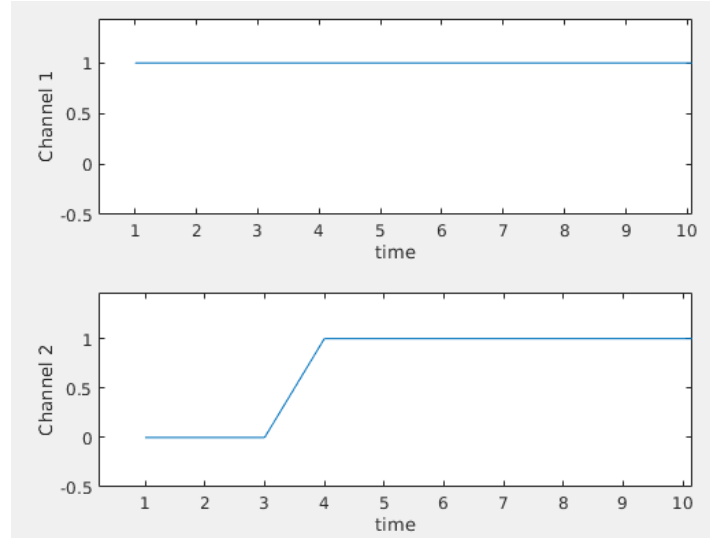
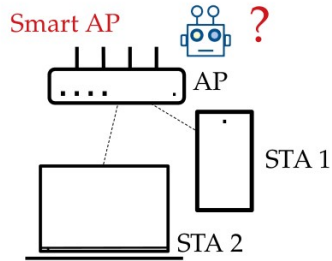
-----



# Case I

channels 1 2

$\alpha=0.2$   
 $\gamma=0.9$



```

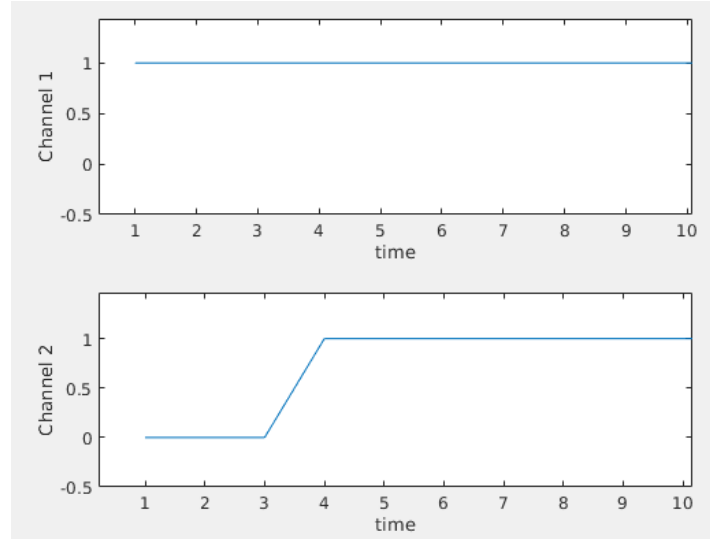
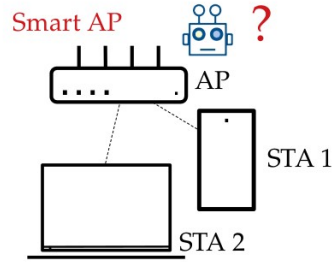
-----Q-table-----
0.0000e+000    0.0000e+000
-200.0000e-003  0.0000e+000
0.0000e+000    0.0000e+000
0.0000e+000    0.0000e+000
-----
Iteration | Current State
2.0000e+000  2.0000e+000
Explore? | Action (1 = stay | 2 = switch)
0.0000e+000  2.0000e+000
Next state | Reward (1 = well done | -1 = :( )
3.0000e+000  1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
800.0000e-003    0.0000e+000    200.0000e-003    1.0000e+000    900.0000e-003    0.0000e+000
-----Q-table-----
0.0000e+000    0.0000e+000
-200.0000e-003  200.0000e-003
0.0000e+000    0.0000e+000
0.0000e+000    0.0000e+000
-----

```

# Case I

channels 1 2

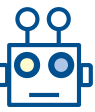
$\alpha=0.2$   
 $\gamma=0.9$

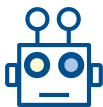


```

-----Q-table-----
  0.0000e+000    0.0000e+000
 -200.0000e-003  200.0000e-003
  0.0000e+000    0.0000e+000
  0.0000e+000    0.0000e+000
-----
Iteration | Current State
  3.0000e+000    3.0000e+000
Explore? | Action (1 = stay | 2 = switch)
  0.0000e+000    1.0000e+000
Next state | Reward (1 = well done | -1 = :( )
  4.0000e+000    -1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
  800.0000e-003    0.0000e+000    200.0000e-003    -1.0000e+000    900.0000e-003    0.0000e+000
-----Q-table-----
  0.0000e+000    0.0000e+000
 -200.0000e-003  200.0000e-003
 -200.0000e-003  0.0000e+000
  0.0000e+000    0.0000e+000
-----

```

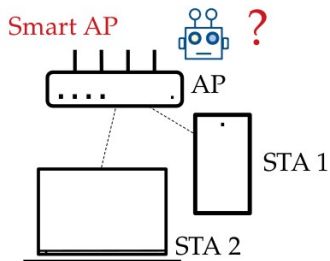




# Case I

channels 1 2

$\alpha=0.2$   
 $\gamma=0.9$



```
-----Q-table-----
0.0000e+000    0.0000e+000
-200.0000e-003  200.0000e-003
-200.0000e-003  0.0000e+000
0.0000e+000    0.0000e+000
```

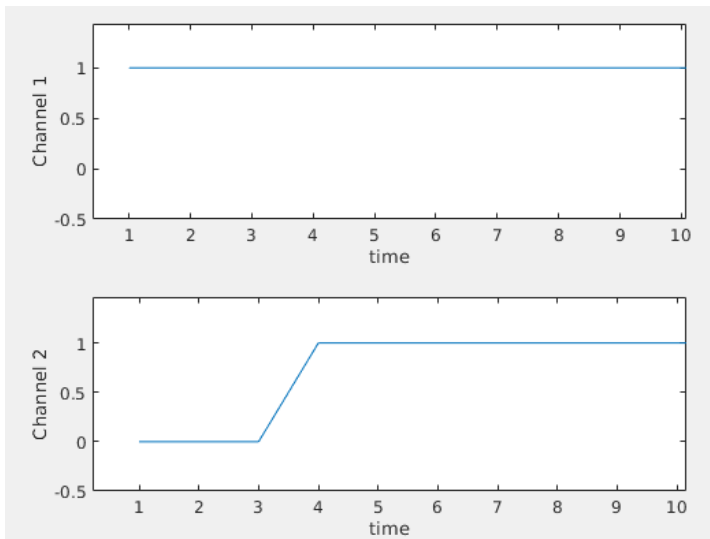
```
-----
Iteration | Current State
4.0000e+000  4.0000e+000
Explore? | Action (1 = stay | 2 = switch)
0.0000e+000  1.0000e+000
Next state | Reward (1 = well done | -1 = :( )
4.0000e+000  -1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
800.0000e-003    0.0000e+000    200.0000e-003    -1.0000e+000    900.0000e-003    0.0000e+000
```

```
-----Q-table-----
0.0000e+000    0.0000e+000
-200.0000e-003  200.0000e-003
-200.0000e-003  0.0000e+000
-200.0000e-003  0.0000e+000
```

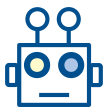
```
-----
Iteration | Current State
5.0000e+000  4.0000e+000
Explore? | Action (1 = stay | 2 = switch)
0.0000e+000  2.0000e+000
Next state | Reward (1 = well done | -1 = :( )
2.0000e+000  -1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
800.0000e-003    0.0000e+000    200.0000e-003    -1.0000e+000    900.0000e-003    200.0000e-003
```

```
-----Q-table-----
0.0000e+000    0.0000e+000
-200.0000e-003  200.0000e-003
-200.0000e-003  0.0000e+000
-200.0000e-003  -164.0000e-003
```

-----



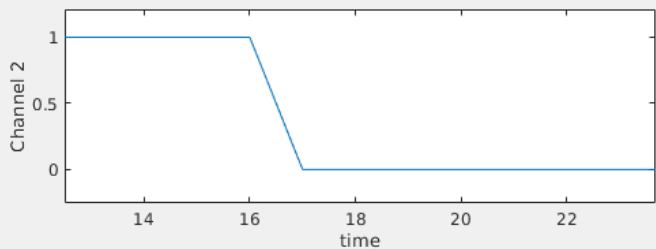
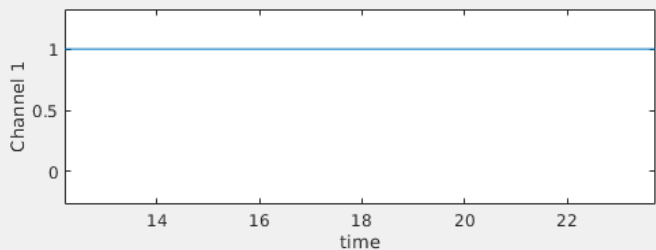
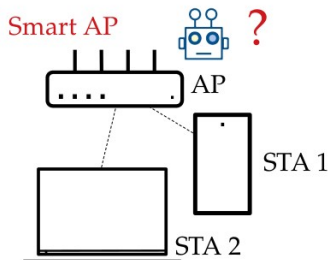




# Case I



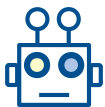
$\alpha=0.2$   
 $\gamma=0.9$



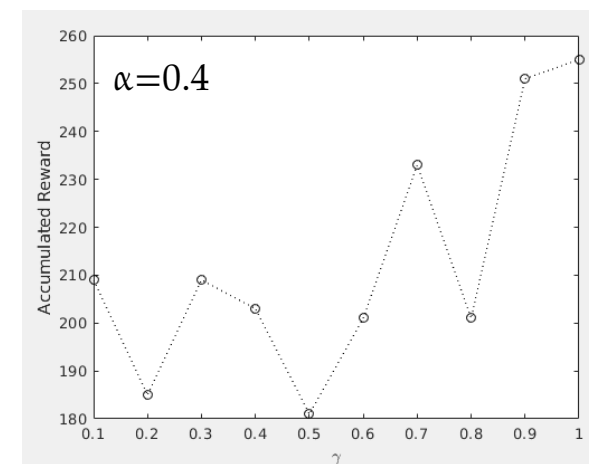
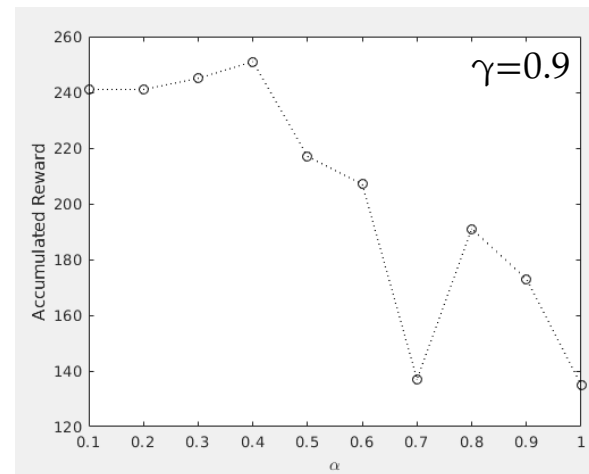
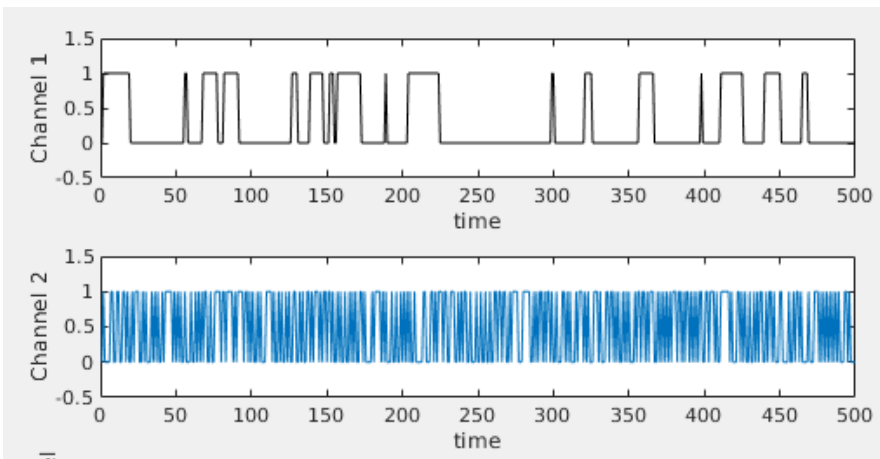
```

-----Q-table-----
  0.0000e+000   0.0000e+000
-709.9346e-003 -219.2809e-003
-200.0000e-003 -239.4706e-003
-581.7296e-003 -680.0567e-003
-----
Iteration | Current State
 19.0000e+000   2.0000e+000
Explore? | Action (1 = stay | 2 = switch)
  0.0000e+000   2.0000e+000
Next state | Reward (1 = well done | -1 = :( )
  3.0000e+000   1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
 800.0000e-003 -219.2809e-003  200.0000e-003   1.0000e+000  900.0000e-003 -200.0000e-003
-----Q-table-----
  0.0000e+000   0.0000e+000
-709.9346e-003 -11.4247e-003
-200.0000e-003 -239.4706e-003
-581.7296e-003 -680.0567e-003
-----
Iteration | Current State
 20.0000e+000   3.0000e+000
Explore? | Action (1 = stay | 2 = switch)
  0.0000e+000   1.0000e+000
Next state | Reward (1 = well done | -1 = :( )
  3.0000e+000   1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
 800.0000e-003 -200.0000e-003  200.0000e-003   1.0000e+000  900.0000e-003 -200.0000e-003
-----Q-table-----
  0.0000e+000   0.0000e+000
-709.9346e-003 -11.4247e-003
  4.0000e-003 -239.4706e-003
-581.7296e-003 -680.0567e-003
-----
Iteration | Current State
 21.0000e+000   3.0000e+000
Explore? | Action (1 = stay | 2 = switch)
  0.0000e+000   1.0000e+000
Next state | Reward (1 = well done | -1 = :( )
  3.0000e+000   1.0000e+000
(1-alfa) * Q_table(state(t),action(t)) + alfa * ( reward(t) + gamma * max(Q_table(state(t+1),:))) =
 800.0000e-003   4.0000e-003  200.0000e-003   1.0000e+000  900.0000e-003   4.0000e-003
-----Q-table-----
  0.0000e+000   0.0000e+000
-709.9346e-003 -11.4247e-003
 203.9200e-003 -239.4706e-003
-581.7296e-003 -680.0567e-003
-----

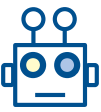
```



# Case II - Test of $\alpha$ and $\gamma$



$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left( \underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)}_{\text{temporal difference}} \underbrace{\phantom{\left( \dots \right)}}_{\text{new value (temporal difference target)}}$$



# Exercise: Congestion Control Video Server

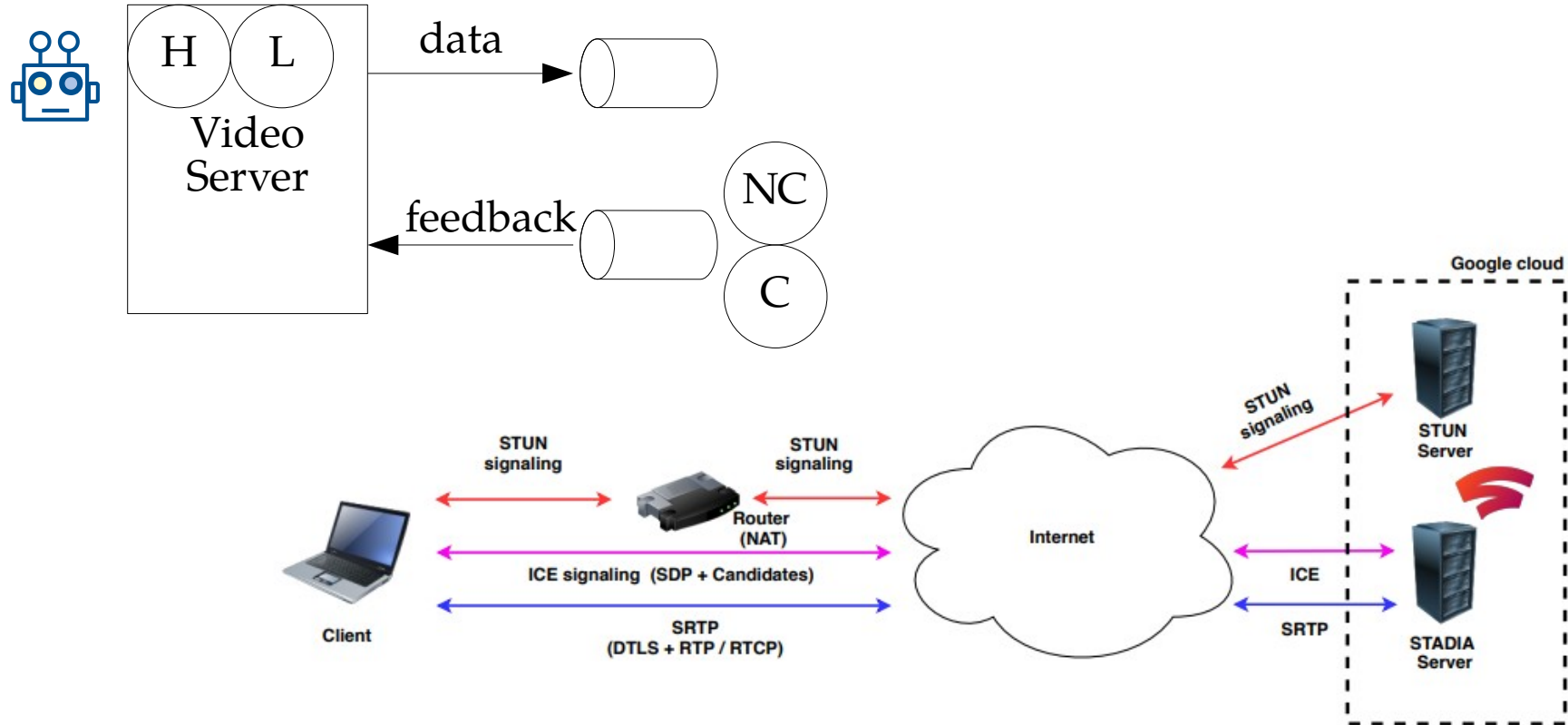
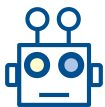


Fig. 1: Google's Stadia: Main components and data streams.



# Google Stadia: WebRTC

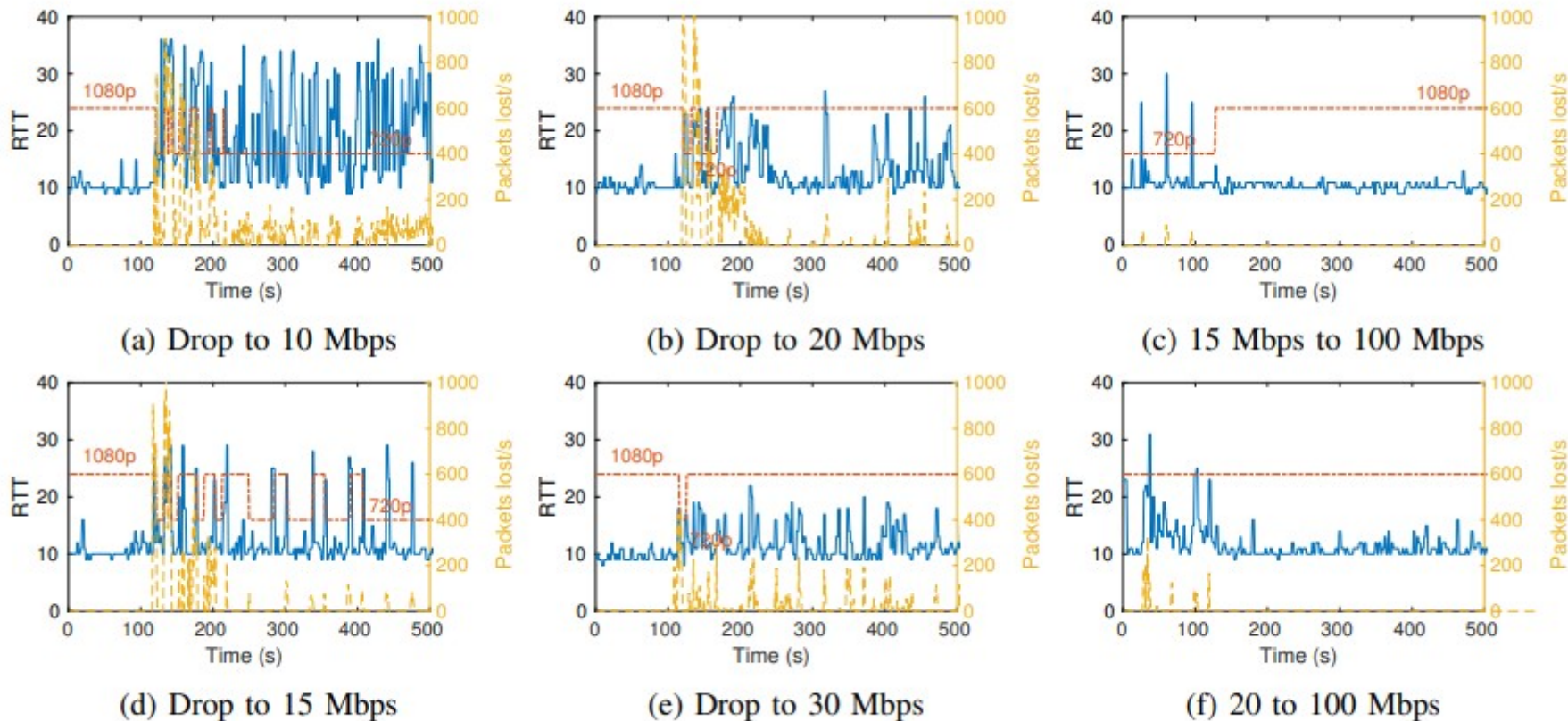
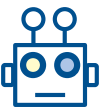
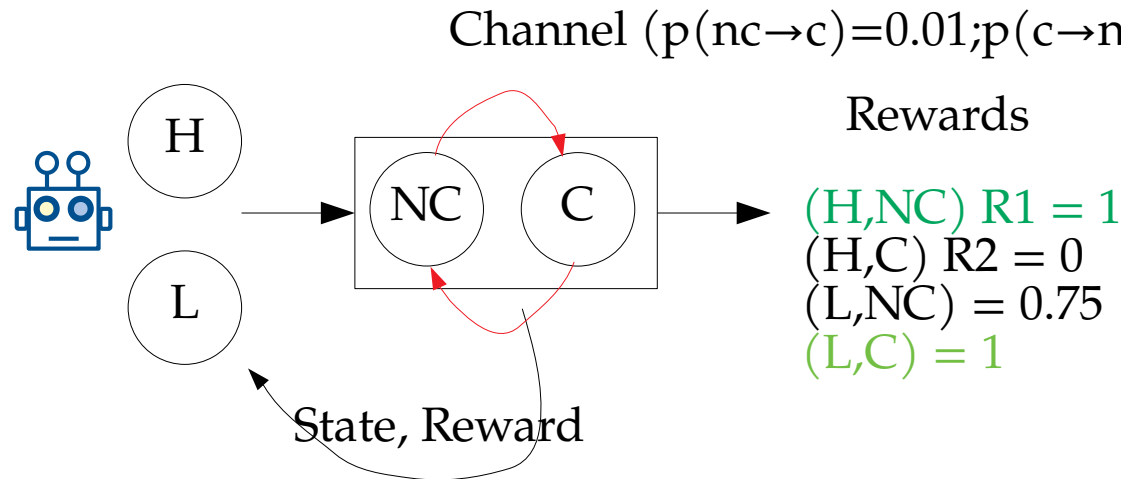


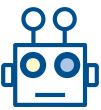
Fig. 11: Round Trip Time (continuous line), video packets lost (dashed line) and resolution (dash-dotted line).



# Exercise: Congestion Control Video Server

- Implement the described scenario, try different channel transitions probabilities
  - Video sender: H (high quality), L (low quality), Channel: NC (Non congested), C (congested)
- Test assuming the agent makes random decisions at each iteration.
- Implement Q-learning, and evaluate if there is any gain. Fine tune Q-learning.





# Activity

- Investigate what is the effect of the rewards (change it, trying to make them consistent ...)
- Will the agent learn a different strategy if  $(H,NC) R1 = 1$   
 $(H,C) R2 = 0$   
 $(L,NC) = 0.25$   
 $(L,C) = 0.75$