

# Let's Play Congestion Control\*

Extended Abstract<sup>†</sup>

Maximilian Bachl  
Institute of Telecommunications  
Vienna, Austria  
maximilian.bachl@tuwien.ac.at

Tanja Zseby  
Institute of Telecommunications  
Vienna, Austria  
tanja.zseby@tuwien.ac.at

Joachim Fabini  
Institute of Telecommunications  
Vienna, Austria  
joachim.fabini@tuwien.ac.at

## ABSTRACT

This paper provides a sample of a  $\LaTeX$  document which conforms, somewhat loosely, to the formatting guidelines for ACM SIG Proceedings.<sup>1</sup>

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; **Redundancy**; Robotics; • **Networks** → Network reliability;

## KEYWORDS

Put some comma-separated keywords here

### ACM Reference Format:

Maximilian Bachl, Tanja Zseby, and Joachim Fabini. 2017. Let's Play Congestion Control: Extended Abstract. In *Proceedings of ACM Woodstock conference (WOODSTOCK'97)*. ACM, New York, NY, USA, 4 pages. [https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

## 1 PARTIAL ACTION ACTOR CRITIC LEARNING

Partial Action Actor Critic Learning (PAL) is a modification of the Asynchronous Advantage Actor Critic algorithm proposed by Mnih et al. [1]. It consists of two Artificial Neural Networks (ANNs), the Actor Network and the Value Network (the Critic part in the abbreviation stands for the Value Network). Given a state, the Actor Networks outputs what it deems to be the optimum action to perform in that certain state. The Value Network estimates what long-term reward can be expected in this state. So an action is considered good if it achieved a long-term reward that is higher than the long-term reward expected by the Value Network and it is considered bad if the reward was lower than expected. The long-term reward is implemented as a moving average of future rewards. So if a high reward can be achieved right now this is more favorable than if it can be achieved in the future. However it can also be beneficial to get a low reward now and instead get a very large one in the future.

\*Produces the permission block, and copyright information

<sup>†</sup>The full version of the author's guide is available as `acmart.pdf` document

<sup>1</sup>This is an abstract footnote

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).  
WOODSTOCK'97, July 1997, El Paso, Texas USA  
© 2017 Copyright held by the owner/author(s).  
ACM ISBN 123-4567-24-567/08/06...\$15.00  
[https://doi.org/10.475/123\\_4](https://doi.org/10.475/123_4)

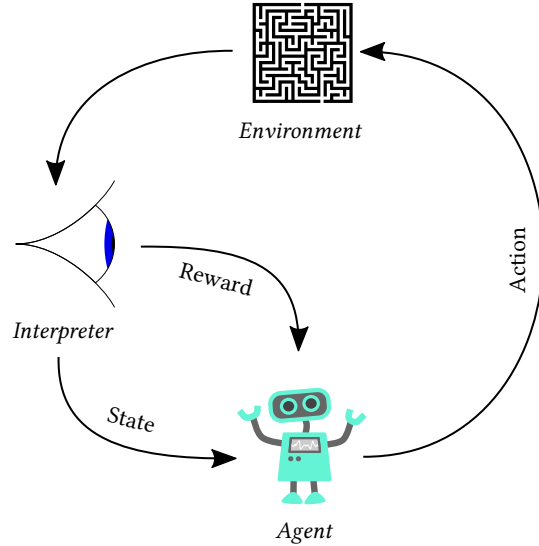


Figure 1: The classical reinforcement learning approach.<sup>a</sup>

<sup>a</sup>adapted from [https://upload.wikimedia.org/wikipedia/commons/1/1b/Reinforcement\\_learning\\_diagram.svg](https://upload.wikimedia.org/wikipedia/commons/1/1b/Reinforcement_learning_diagram.svg)

### 1.1 Algorithm

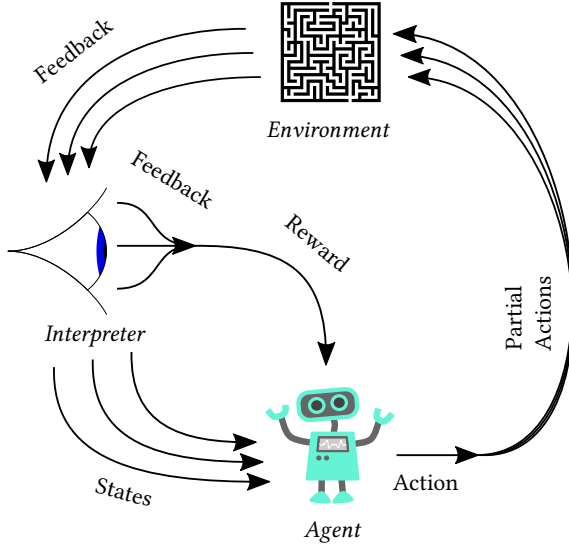
In Algorithm 1 we show the code that runs in one of the learning threads. As proposed in literature, we use several concurrent threads to improve training performance. A major difference to previous approaches to reinforcement learning is the fact that we allow more actions to be taken before the previous actions' rewards have been received. In classical reinforcement learning, an action is always followed by a reward and a reward is always followed by an action. In our proposed concept, however, it is possible to take new actions while the previous action hasn't received its reward yet.

Another major difference in PAL is that one action generates a number of partial actions ( $\geq 0$ ) (see Figure 2). Each partial action generates feedback upon interacting with the environment. Upon receiving feedback for a partial action, the agent determines the current state and triggers a new action. When all partial rewards of one action were received, the agent combines them to form the reward and updates the value and actor networks (see Algorithm 2).

### 1.2 Value Network

The value network estimates the expected long-term reward given a state  $s_t$ .

Let  $r_t$  be the reward that was received at time  $t$ .  $V(s_t; \theta_v)$  designates the expected mean reward in state  $s_t$  given the parameters



**Figure 2: Partial Action Actor Critic Learning: An action consists of zero or more partial actions which trigger partial rewards upon interacting with the environment. Each partial reward updates the state. The value and actor networks are updated upon receiving all partial rewards of one action.<sup>a</sup>**

<sup>a</sup>adapted from [https://upload.wikimedia.org/wikipedia/commons/1/1b/Reinforcement\\_learning\\_diagram.svg](https://upload.wikimedia.org/wikipedia/commons/1/1b/Reinforcement_learning_diagram.svg)

(neural network weights) of the value network  $\theta_v$ . Then, with  $\gamma$  being the roll-off factor, which designates the influence that future reward has on the moving average, (set to 0.99), we define the average reward as

$$R_t = \left( \sum_{i=0}^{k-1} \gamma^i r_{t+i} \right) + \gamma^k V(s_{t+k}; \theta_v),$$

where  $k$  is upper-bounded by  $t_{\max}$  ( $t_{\max}$  is a fixed hyperparameter that indicates how many steps should be performed before updating the neural network). So  $R$  is simply a moving average of rewards. However, usual moving averages take into account values from the past while this one uses values from the future; it runs reversely.

The loss function, which the value network tries to minimize, is the square of the difference of the actual average reward received and the expected average reward

$$l_{v,t} = (R_t - V(s_t; \theta_v))^2.$$

### 1.3 Actor Network

Let  $a_t$  be the action that was taken at time step  $t$ ,  $v_t$  the value that the value network estimated as the future reward given the current state  $s_t$  at time step  $t$  ( $v_t$  is just an abbreviation for  $V(s_t; \theta_v)$ ) and  $\theta_a$  the parameters (neural network weights) of the actor network. Furthermore, let  $\beta$  be a factor that specifies the importance of the entropy  $H$ .  $\pi$  designates the probability density function, which means that  $\pi(a_t | s_t; \theta_a)$  is the value of the probability density function of taking action  $a_t$  in state  $s_t$  with the current weights of

**Algorithm 1** Partial Action Actor Critic Learning – pseudocode for each learning thread. Only  $\theta_{a,g}$  and  $\theta_{v,g}$  are shared between the threads. They are initialized randomly in the beginning. **TODO: Update algorithm**

```

1: loop
2:    $l_{\text{actions}} \leftarrow []$ 
3:    $l_{\text{states}} \leftarrow []$ 
4:    $l_{\text{values}} \leftarrow []$ 
5:    $l_{\text{rewards}} \leftarrow []$ 
6:    $l_{\text{estimatedValues}} \leftarrow []$ 
7:    $l_{\text{snapshots}} \leftarrow []$ 
8:    $t \leftarrow 0$ 
9:    $s_0 \leftarrow \text{initialState}()$ 
10:   $\theta_a \leftarrow \theta_{a,g}$ 
11:   $\theta_v \leftarrow \theta_{v,g}$ 
12:  repeat
13:     $l_{\text{states}}.\text{append}(s_t)$ 
14:     $l_{\text{estimatedValues}}.\text{append}(V(s_t; \theta_v))$ 
15:    if  $t \bmod t_{\max} = 0$  then
16:       $\theta_a \leftarrow \theta_{a,g}$ 
17:       $\theta_v \leftarrow \theta_{v,g}$ 
18:       $l_{\text{snapshots}}.\text{append}((\theta_a, \theta_v))$ 
19:    end if
20:    Sample  $a_t$  from the
      Actor Network's probability distribution
21:     $l_{\text{actions}}.\text{append}(a_t)$ 
22:     $l_{\text{partialActions},t} \leftarrow \text{partialActions}(a_t)$ 
23:    for all partial actions  $a_{p,n,t}$  in  $l_{\text{partialActions},t}$  do
24:      Take partial action  $a_{p,n,t}$ 
25:    end for
26:    Receive feedback  $r_{p,n,t'}$ 
      where  $t' \leq t$  and  $0 \leq n \leq \#(l_{\text{partialActions},t'})$ 
27:    if all feedback of  $a_{t'}$  was received then
28:       $r_{t'} \leftarrow \text{reward w.r.t all feedback } r_{p,n,t'}$ 
29:       $l_{\text{rewards}}.\text{append}(r_{t'})$ 
30:      if  $\#(l_{\text{rewards}}) > t_{\max}$  or the episode is over then
31:        COMPUTE GRADIENTS
32:      end if
33:    end if
34:    Generate  $s_{t+1}$  using  $r_{p,t'}$ 
35:     $t \leftarrow t + 1$ 
36:  until reaching the end the episode
37: end loop

```

the actor network  $\theta_a$ . Then

$$l_{a,t} = -\log(\pi(a_t | s_t; \theta_a))(R_t - v_t) - \beta H(\pi(s_t; \theta_a))$$

is the actor loss that we try to minimize.

## 2 CONGESTION CONTROL SPECIFICS

The classical Reinforcement Learning assumes that a reward follows an action and vice-versa (see Figure 1). However, in the case of congestion control, it is desirable to perform a new action without having received a reward for the previous action. For example,

**Algorithm 2** Partial Action Actor Critic Learning – procedure which computes and applies the gradients.

```

1: function COMPUTEGRADIENTS
2:    $t_{\text{end}} = \min(t_{\text{max}}, \#(l_{\text{rewards}}))$ 
3:    $\theta_{\text{backup}} \leftarrow (\theta_a, \theta_v)$ 
4:    $\theta_a, \theta_v \leftarrow l_{\text{snapshots}}[0]$ 
5:    $R_{i+1} \leftarrow \frac{l_{\text{estimatedValues}}[t_{\text{end}}]}{1-\gamma}$ 
6:   for  $i \leftarrow t_{\text{end}} - 1, 0$  do
7:      $R_i \leftarrow r_i + \gamma R_{i+1}$ 
8:      $a \leftarrow l_{\text{actions}}[i]$ 
9:      $s \leftarrow l_{\text{states}}[i]$ 
10:     $v \leftarrow l_{\text{values}}[i]$ 
11:     $d\theta_v \leftarrow \frac{\partial(R_i - v)}{\partial \theta_v}$ 
12:     $d\theta_v \leftarrow - \frac{\frac{\partial \log(\pi(a | s; \theta_a))(R_i - v)}{\partial \theta_v}}{\frac{\partial \beta H(\pi(s; \theta_a))}{\partial \theta_a}}$ 
13:     $\theta_{a,g} \leftarrow \theta_{a,g} + d\theta_a$ 
14:     $\theta_{v,g} \leftarrow \theta_{v,g} + d\theta_v$ 
15:  end for
16:  Remove first  $t_{\text{end}}$  elements from
     $l_{\text{actions}}, l_{\text{states}}, l_{\text{values}}, l_{\text{rewards}}, l_{\text{estimatedValues}}$ 
17:  Remove first element from  $l_{\text{snapshots}}$ 
18:   $\theta_a, \theta_v = \theta_{\text{backup}}$ 
19: end function

```

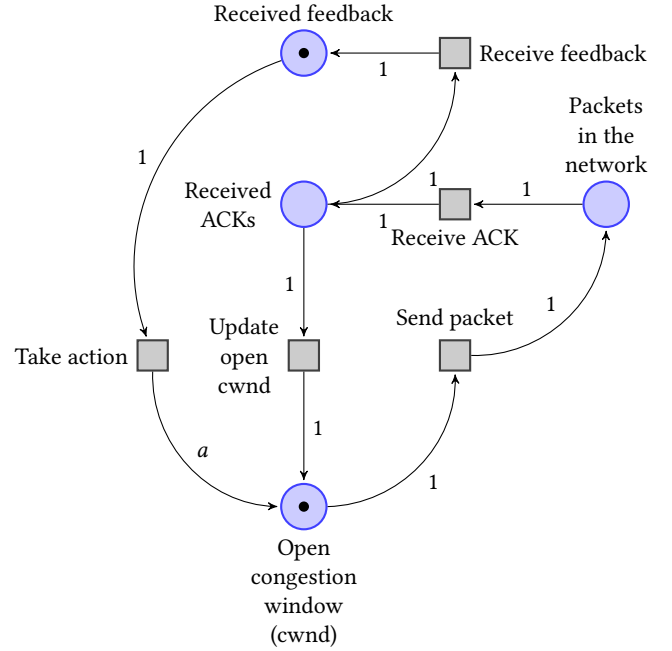
imagine that you receive two acknowledgements directly after each other. For each of these two acknowledgements an action has to be performed but it does not seem feasible for the second action to wait until the first action has received a reward. The Asynchronous Actor Critic framework as described by [1] cannot be applied to congestion control as it assumes that actions and rewards are synchronized and thus we use the proposed Partial Action Actor Critic Learning (see section 1). We describe the concept of using reinforcement learning for congestion control with a petri net (see Figure 3).

In the following a time step  $t$  corresponds to the reception of an acknowledgement. The beginning of the flow, before any packet is sent, corresponds to time step 0.

- $s_t$  The state describes the current “congestion state”. Various features that indicate congestion are included in it. Each time an acknowledgement is received, the state is updated and the actor network is asked for the next action.
- $a_t$  Based on a given state and the history of previous states, the actor network returns an action  $a_t$ , which is a real number ( $\geq 1$ ) that stands for the congestion window to be used until the next acknowledgement is received.
- $r_t$  The reward is a tuple of at least one reward metric. These are discussed in more detail in subsection 2.1
- $v_t$  The value is a tuple of the expected average reward estimated by the value network (see subsection 2.1) for each of the reward metrics.

We actually use four different types of reward:

- $r_{\text{packet}, t}$  is the number of the packets that the sender sent during time step  $t$  and that were not lost (so they were acknowledged at some point by the receiver).



**Figure 3:** A petri net describing the congestion control mechanism. We start with one token in the state *Open congestion window*, which means that initially we can send one packet. We start with another token in the state *Received feedback*, which means that we can take an action right in the beginning. The concept is the following: If there is at least one token in the *Open congestion window*, send a packet. The acknowledgement (or timeout if the packet gets lost) for the packet is received by the *Open congestion window*, meaning that when we get an ACK we can send another packet (The *Open congestion window* signifies the congestion window minus the packets that are currently unacknowledged.). Furthermore the state *Received feedback* receives each ACK/timeout and takes an action. *Take action* adds  $a$  tokens to the *Open congestion window*, where  $a$  is a real number (possibly also negative); thus in each state there can also be a real number of tokens (e.g. 2.34 tokens are possible in the *Open congestion window*). However, we define that the congestion window can never be smaller than 1. **TODO: Draw as 3 times 3 square**

$r_{\text{byte}, t}$  is the sum of the bytes of the packets that the sender sent and that were not lost. Example: During  $t$ , three packets were sent with 300, 200, and 1500 bytes each so  $r_{\text{byte}, t} = 2000$

$r_{\text{delay}, t}$  is the sum of the round trip times of the packets that the sender sent and that were not lost.

$r_{\text{duration}, t}$  is the sum of the time between receiving the last packet and receiving this packet (“inter-receive time”) for the packets that the sender sent and that were not lost.

## 2.1 Value Network

In the case of congestion control, the loss function  $l_{v,t}$  of the value network is actually the sum of the square of the difference for each of the four moving averages for each type of reward:

$$\begin{aligned} l_{v,t} = & \left( R_{\text{packet},t} - V_{\text{packet}}(s_t; \theta_v) \right)^2 \\ & + \left( R_{\text{byte},t} - V_{\text{byte}}(s_t; \theta_v) \right)^2 \\ & + \left( R_{\text{delay},t} - V_{\text{delay}}(s_t; \theta_v) \right)^2 \\ & + \left( R_{\text{duration},t} - V_{\text{duration}}(s_t; \theta_v) \right)^2 \end{aligned}$$

## 2.2 Actor Network

The actor network outputs two parameters: The mean of a normal distribution  $\mu$  and its standard deviation  $\sigma$ .

Each time an action  $a_t$  is requested, a value  $X$  is sampled from the current normal distribution defined by the parameters  $\mu$  and  $\sigma$ :

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

$a_t$  is the sampled value  $X$  at time  $t$  and the window is incremented by  $a_t$ ; however, the window can never be smaller than 1. At the beginning of a flow the window is 1 as well.

With  $H$  being the entropy, the actor network minimizes the loss

$$\begin{aligned} l_{a,t} = & -\log(\pi(a_t | s_t; \theta_a)) \\ & \left( \log\left(\frac{R_{\text{byte},t}}{v_{\text{byte},t}}\right) - \log\left(\frac{v_{\text{byte},t}}{v_{\text{duration},t}}\right) \right) - \left( \frac{R_{\text{delay},t}}{R_{\text{packet},t}} - \frac{v_{\text{delay},t}}{v_{\text{packet},t}} \right) \\ & + \beta H(\pi(s_t; \theta_a)). \end{aligned}$$

## 2.3 Example

- (1) Receive an ACK
- (2) Update the internal state to reflect the information that this ACK provides (for example, the state consists of the round trip time experienced by this packet, the time since we receive the last ACK etc.)
- (3) Sample an increase for the congestion window from the Actor Network (e.g. 0.24) and add it to the current congestion window.
- (4)(a) If this was the last ACK, which completes a previous action (e.g. a previous action resulted in 3 packets being sent; so if this is the ACK for the third (and last) packet of this previous action we got all the ACKs to compute the reward.)
- (b) If we now have  $t_{\text{max}}$  rewards (set to 20) then we also update the neural network.

## 2.4 Convergence

We use

$$\text{Reward} = \log(\text{throughput}) - \delta \text{delay}$$

as our metric of reward, where  $\delta$  is a parameter that allows us to change the importance of the throughput vs. the delay. We take the logarithm of the throughput to ensure that flows which have small throughput have a higher incentive to increase it as flows which

already have a large one (because the derivative of the logarithm decreases with increasing value). We take minus the delay to punish inducing delay. We do not take the logarithm because we think that it is equally bad if a sender behind a satellite connection with a minimum RTT of 500 ms adds 1 ms of delay as if a sender connected using Ethernet with a minimum RTT of 10 ms adds 1 ms of delay.

There are two conditions that we have to take into account to ensure convergence:

(1)

$$\left( \frac{\partial}{\partial w} \text{throughput} \right) (1) > \left( \frac{\partial}{\partial w} \delta \text{delay} \right) (1)$$

(2)

$$\frac{\partial}{\partial w} (\delta)$$

TODO: Finish second constraint for convergence.

TODO: Add figure to visually explain that

## 3 RESULTS

We train on the following scenario:

Parameter	Value	Distribution
Two-way propagation delay	100 ms	constant
Bottleneck bandwidth	3 Mbit/s	constant
Number of senders	2	constant
Flow length	$\infty$	constant
Simulation duration	50 s–150 s	uniform
Buffer size	$\infty$	constant
Stochastic loss prob.	0%	constant

**Table 1: The parameters of our very first simple experiment**

TODO: Insert results when comparing to regular TCP

## REFERENCES

- [1] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*. 1928–1937.