# Deep Reinforcement Learning Empowered Rate Selection of XP-HARQ

Da Wu, Jiahui Feng, Zheng Shi, Hongjiang Lei, Guanghua Yang, and Shaodan Ma

*Abstract*—The complex transmission mechanism of cross-packet hybrid automatic repeat request (XP-HARQ) hinders its optimal system design. To overcome this difficulty, this letter attempts to use the deep reinforcement learning (DRL) to solve the rate selection problem of XP-HARQ over correlated fading channels. In particular, the long term average throughput (LTAT) is maximized by properly choosing the incremental information rate for each HARQ round on the basis of the outdated channel state information (CSI) available at the transmitter. The rate selection problem is first converted into a Markov decision process (MDP), which is then solved by capitalizing on the deep deterministic policy gradient (DDPG) algorithm with the prioritized experience replay. The simulation results finally corroborate the superiority of the proposed XP-HARQ scheme over the conventional HARQ with incremental redundancy (HARQ-IR) and the XP-HARQ with only statistical CSI.

*Index Terms*—Cross-packet hybrid automatic repeat request (XP-HARQ), deep reinforcement learning (DRL), outdated channel state information, rate selection.

## I. INTRODUCTION

Hybrid automatic repeat request (HARQ) is one of the key technologies that is capable of offering reliable transmissions. However, this benefit is essentially reaped at the price of large transmission delay, which is unfavorable for fulfilling the ultra-reliable and low-latency communications (URLLC). To resolve such a dilemma, there is an urgent need to develop a flexible HARQ transmission mechanism that could be reconfigurable to meet diverse URLLC requirements. In this letter, we focus on the cross-packet HARQ (XP-HARQ) that is an evolutionary version of HARQ with high spectral efficiency, albeit at the price of high complexity [1], [2], [3]. Unlike the conventional HARQ schemes, new information bits are introduced in retransmissions such that

of its high complexity [4], [5]. Unlike the transmission of HARQ, each message gives information to the delivery of another message, especially in such a case that all prior wireless resource substantially exploited by HARQ is boosted while the average end-to-sink rate boosted. Recently, the XP-HARQ has been investigated for its high throughput and ultra-low latency. As a consequence, the spectral efficiency of HARQ is boosted in XP-HARQ such as average transmission delay is reduced.

Recently, to investigate the throughput and reliability of XP-HARQ, efforts have been made to accurately evaluate and optimally design XP-HARQ. For instance, the authors in [6] examined the long term average throughput (LTAT) of XP-HARQ with puncturing throughput optimized with an XP-HARQ was verified. The adaptive modulation coding scheme was introduced to implement XP-HARQ to attain XP-HARQ in [7] to offer with capacity same XP-HARQ were analyzed for buffer-aided XP-HARQ. However the per-round rate of XP-HARQ was optimized [8] was obtained by conducting Monte-Carlo simulations and lacked insightful analysis. To further this vast the optimal design of XP-HARQ, from an information theoretic probability channel derived in closed-form for XP-HARQ over the independent Rayleigh fading channels in [8], with full diversity of XP-HARQ was provided. However, the result should be simple that only model, works out attempted to devise the completion of HARQ assistance the optimal design of XP-HARQ, but to identify driving scheme complicated fading channels designed to mitigate the above situation. As for to HARQ systems, deep reinforcement learning (DRL) from DDPG algorithm of XP-HARQ to maximize the fading channel utilizing the be noticed that only a few works attempted to devise the DRL into the HARQ. As such HARQ using the DRL method. Particularly, in [9] an DRL enabled LTAT scheduling policy was designed to considering the age of information (AoI) of HARQ systems. In [9] deep deterministic policy gradient (DDPG) algorithm was used to minimize the the throughput maximizing the incremental redundancy bits. solved by namely DDPG extension of the DRL method. By generate HARQ, the ideal is next to be proposed. This XP-HARQ maximizes the LTAT by a superior rate selection of HARQ with outdated channel state (HARQ-IR) in CSI XP-

Fig. 1. An example of a XP-HARQ scheme with $K = 3$.

## II. System Model

### A. XP-HARQ

This letter considers a point-to-point communication system, in which XP-HARQ is adopted to enable the retransmissions of the message. To start, this section delineates the system model, including the XP-HARQ transmission mechanism, the channel model, performance metrics, and the rate selection problem. For notational simplicity, let $n(t) \in \mathbb{Z}^+$ and $\kappa(t) \in [1, K]$ be the functions that map the time slot $t$ to the current HARQ cycle and the current transmission round, respectively. In the initial transmission round of the $n(t)$-th HARQ cycle, the message $m_{n(t),1}$ is encoded as a codeword $\mathbf{x}_{n(t),1}$ with the transmission rate $R_1$. The received signal $\mathbf{y}_{n(t),1}$ reads as

$$\mathbf{y}_{n(t),1} = \sqrt{P_1} h_{n(t),1} \mathbf{x}_{n(t),1} + \mathbf{n}_{n(t),1}, \quad (1)$$

where $h_{n(t),1}$ denotes the channel coefficient of the $n(t)$-th HARQ cycle with transmission round $\kappa(t)$. According to the coding strategy of XP-HARQ, if a positive acknowledgement (ACK) will be sent back to confirm the successful reception of $m_{n(t),1}$ and let the next HARQ cycle with index $t_{n(t)+1}$ will be triggered immediately. Otherwise, a negative acknowledgement (NACK) will be fed back to initiate the retransmission. According to the coding strategy of XP-HARQ, where only the incremental information bits are retransmitted, a new information-bits are introduced in the $\kappa$-th transmission. The XP-HARQ signal is substantially exploited wireless resources. Accordingly, in the $\kappa(t)$-th transmission of the $n(t)$-th HARQ cycle, the previously failed messages $m_{n(t),1}, \cdots, m_{n(t),\kappa(t)-1}$ are combined with the currently received message $m_{n(t),\kappa(t)}$ to form a longer message $m_{n(t),[\kappa(t)]}$. The concatenated message $m_{n(t),[\kappa(t)]}$ is encoded as a codeword $\mathbf{x}_{n(t),\kappa(t)}$ with a nominal transmission rate $\sum_{\kappa=1}^{\kappa(t)} R_\kappa$ and $P_1$, respectively, which increment over the transmission. The messages $R_{\kappa(t)}$ originates from the jointly information bits involved in the $\kappa$-th transmission. Therefore, the signal $\mathbf{y}_{\kappa(t)}$ stops received in the $\kappa(t)$-th round once the receiver XP-HARQ cycle

$$\mathbf{y}_{n(t),\kappa(t)} = \sqrt{P_{\kappa(t)}} h_{n(t),\kappa(t)} \mathbf{x}_{n(t),\kappa(t)} + \mathbf{n}_{n(t),\kappa(t)}, \quad (2)$$

where $h_{n(t),\kappa(t)}$, $\mathbf{n}_{n(t),\kappa(t)}$, and $P_{\kappa(t)}$ follow the similar definitions as $h_{n(t),1}$, $\mathbf{n}_{n(t),1}$, and $P_1$, respectively, which are omitted here to save space. The messages $m_{n(t),1}, \cdots, m_{n(t),\kappa(t)}$ are jointly decoded by using the observations $y_1, \cdots, y_{\kappa(t)}$. The current XP-HARQ cycle stops and the next process begins once the receiver succeeds in reconstructing all the previously delivered messages or the maximum number of HARQ transmission attempts $\kappa(t)$ is used. Interested readers are referred to [?] for more details of the encoding/decoding implementation of XP-HARQ.

### B. Channel Model

This letter considers time-correlated Rayleigh flat-fading channels, where the channel keeps constant during each codeword transmission slot and changes time-dependently across consecutive transmission slots. We define $t$ as the index of the time slot in the sequel. For notational simplicity, we use the notation $h_t$ to represent $h_{n(t),\kappa(t)}$. As a commonly used time-correlated channel model that takes place in the environment of low-to-medium mobility, $h_t$ is modeled according to a first-order Gauss-Markov process as [?], i.e.,

$$h_t = \rho h_{t-1} + \sqrt{1-\rho^2}\, w_t, \quad (3)$$

where $\rho$ is the correlation coefficient between $h_t$ and $h_{t-1}$. $w_t \sim \mathcal{CN}(0, \sigma^2)$ denotes the channel discrepancy and is independent of $h_{t-1}$. In order to account for the impact of channel aging, the outdated channel state $h_{t-1}$ is sent back to the transmitter.

### C. Performance Metrics

1) Outage Probability: The outage probability is an essential performance metric for evaluating the system reliability. The outage probability of XP-HARQ is the probability that the accumulated mutual information in each HARQ round is below the transmission rate. More specifically, the outage probability of XP-HARQ after $K$ HARQ rounds is given by [?]

information bits. The long-term transmission throughput, i.e., the specific frequently used outage probability of XP-HARQ after $K$ HARQ rounds is given by [?]. The LTAT of XP-HARQ system is defined as [?]

$$f_K = \Pr\left(I_1 < R_1, I_2 < R_2^\Sigma, \cdots, I_K < R_K^\Sigma\right), \quad (4)$$

where $I_k = \sum_{l=1}^{k} \log_2\left(1 + |h_l|^2 P_l/\sigma^2\right)$ stands for the accumulated mutual information until the $l$-th transmission.

*2) Long Term Average Throughput:* The long term average throughput (LTAT) is a frequently used performance metric to evaluate the expected throughput of HARQ systems [?]. The LTAT of XP-HARQ system is defined as [?]

### D. Maximization of LTAT

This paper aims to maximize the LTAT through optimal rate selection if only the aged channel state information (CSI) is available at the transmitter. The optimization problem of the transmission rates can be formulated as

$$\eta_K = \lim_{T\to\infty} \frac{\mathcal{R}(T)}{T} = \frac{\sum_{k=1}^{K} R_k (f_{k-1} - f_K)}{1 + \sum_{k=1}^{K-1} f_k}, \quad (5)$$

where $\mathcal{R}(T)$ refers to the total number of successfully received information bits till time $t$, and the second equality in (??) is derived in [?], [?] by capitalizing on the renewal theory if only the statistical CSI is available at the transmitter.

$$\max_{R_1,\cdots,R_K} \eta_K \quad (6)$$
$$\text{s.t.} \quad 0 \le R_k \le \bar{R}, k \in [1, K],$$

where the transmission rate $\{R_k, k \in [1, K]\}$ is upper bounded by $\bar{R}$ to avoid frequent outages because of the limited resources. However, due to the time correlation among fading channels in (??) and the involved outage definition in (??), it is hardly possible to get the explicit LTAT maximization problem in (??) with the conventional optimization tools. To overcome this difficulty, we recourse to the deep reinforcement learning (DRL) for the optimal solution of the transmission rate.

## III. DRL EMPOWERED RATE SELECTION

Due to the rapid change of time-varying fading channels, it results in a prohibitively high system overhead to acquire the instantaneous CSI. Therefore, we assume that only the outdated and statistical CSIs are available at the transmitter, including the channel state of the previous slot $h_{t-1}$ and the correlation coefficient $\rho$. Moreover, the transmission rate of the current transmission round for XP-HARQ is determined by the transmission status (success or failure), rates, and channel states in the previous transmission rounds. Towards this end, the proposed optimization problem is transformed into a Markov decision process (MDP), which can be solved with DRL methods.

### A. Problem Reformulation for MDP

By using the definition of the LTAT and replacing the limit operation with the expectation (the time average converges to the ensemble average for ergodic processes), the original problem (??) can be reformulated as

$$\max_{R(t),\,n(t)} \mathbb{E}\left(\frac{\mathcal{R}(T)}{T}\right) \quad (7)$$
$$\text{s.t.} \quad 0 \le R(t) \le \bar{R},$$

With the problem reformulation of (??), the adaptive rate selection scheme can be modeled as an MDP, which makes the sequential decision. But the effective transmission rate $\mathcal{R}_{n(t),\kappa(t)}$ denotes the effective transmission rate for the successfully received information bits after step $t$ to the process in the $n(t)$-th HARQ cycle. According to the Shannon theory, the successful decoding occurs if and only if the transmission rate is less than the channel capacity. Therefore, $\mathcal{R}_{n(t),\kappa(t)}$ is received outdated environment $\mathcal{E}$. By mapping the optimal rate selection of XP-HARQ as an MDP, the states, actions, and rewards are designed as follows.

$$\mathcal{R}_{n(t),\kappa(t)} = \begin{cases} R_{\kappa(t)}^\Sigma, & I_{\kappa(t)} \ge R_{\kappa(t)}^\Sigma \\ 0, & \text{else} \end{cases} \quad (8)$$

*1) State $s_t$:* To capture the channel aging effect, the historical channel state $h_t$ is considered into the observation of environment. Moreover, the decoding status of XP-HARQ essentially depends on the accumulated mutual information and rate. Accordingly, the state $s_t$ is a vector consisting of the previously accumulated transmission rate and effective information rate of the successfully XP-HARQ and the aged channel HARQ cycle, i.e.,

$$s_t = \begin{cases} \left(R_{\kappa(t-1)}^\Sigma, \mathcal{R}_{n(t),\kappa(1)} = n(t)\right), & n(t-1) = n(t) \\ (0,0), & \text{else} \end{cases} \quad (9)$$

By noticing the continuous space of the state and actions, the MDP problem can be solved with the DRL, which combines the reinforcement learning and deep neural networks to learn the policy. The details are deferred to the next subsection.

*2) Action $a_t$:* The action is defined as the effective transmission rate for the new information bits in the next HARQ round, i.e.,

$$a_t = R(t). \quad (10)$$

*3) Reward $r_t$:* The reward function can be defined as the effective transmission rate of the successfully received information bits for the current HARQ cycle $n(t)$, i.e.,

$$r_t = \mathcal{R}_{n(t),\kappa(t)}. \quad (11)$$

### B. DRL Empowered Rate Selection

A DRL based rate selection scheme is proposed for the LTAT maximization of the XP-HARQ. By considering the continuous state and action spaces, a deep deterministic policy gradient (DDPG) with prioritized experience replay will be applied to develop the rate selection framework, as shown in Fig. ??. This framework consists of four neural networks, i.e., two policy networks (also termed as the actor network)

By noticing the continuous spaces of the evaluation and policy networks, the MDP problem can be solved with the DRL, which combines the reinforcement learning and deep neural networks. The details are referred to the next subsection. To address the overestimation issue, and these neural networks are parameterized by $\theta$, $\theta^-$, $\omega$, and $\omega^-$. In addition, for the stability and fast convergence, a prioritized experience reply memory pool $\mathcal{M}$ is adopted to collect the agent's experience tuple $e_t = (s_i, a_i, r_i, s_{i+1})$, at each time $t$.

## B. DRL Empowered Rate Selection

A DRL-based rate selection scheme is proposed for the KTAT maximization of the XP-HARQ. By considering the continuous state and action spaces, a deep deterministic policy gradient (DDPG) with prioritized experience replay will be applied to develop the rate selection framework, as shown in Fig. (??). This framework consists of four neural networks, i.e., two policy networks (also termed as the actor networks described $\mu(s_t; \theta)$ and $\mu(s_{t+1}; \theta^-)$) and two evaluation networks (also termed as the critic network $Q(s_t, a_t; \omega)$ and $Q(s_{t+1}, a_{t+1}; \omega^-)$), wherein the target-evaluation and target-policy networks are used to calculate the temporal-difference ("TD") target to address the overestimation issue, and these neural networks are parameterized by $\theta$, $\theta^-$, $\omega$, and $\omega^-$. In addition, for the stability and fast convergence, a prioritized experience reply memory pool $\mathcal{M}$ is adopted to collect the agent's experience tuple $e_t = (s_i, a_i, r_i, s_{i+1})$, at each time step, the four neural networks will be updated with a mini-batch of experience samples $\mathcal{B}_t$ that are drawn from $\mathcal{M}$ according to the priority of the playback experience, that is, $e_t \sim \mathcal{P}(\mathcal{M})$ for $\forall e_t \in \mathcal{B}_t$, where $\mathcal{P}$ is the probability function defined in (??). In what follows, priority experience playback mechanism and the training processes of the four neural networks are described in detail.



Fig. 2. The DDPG network for Rate Selection of XP-HARQ

*1) Prioritized Experience Replay:* In contrast with the uniform random experience replay, the prioritized experience replay is capable of accelerating the learning process and enhancing the training stability [?]. According to the prioritized sampling strategy, the sampling probability $p_i$ of the tuple $e_i$ is proportional to the absolute value of TD error $\delta_i$, i.e.,

$$ w_i \propto (|\mathcal{B}_i| p_i) \epsilon^{-\beta}, \tag{12} $$

where $\epsilon$ is a positive constant to avoid a zero probability, $\delta_i = Q(s_i, a_i; \omega) - r_i - \gamma Q(s_{i+1}, a_{i+1}; \omega^-)$ denotes the TD error, and $\gamma$ is the discount factor.

*2) Evaluation Network:* The evaluation network aims to approximate the actual state-action function $Q_\pi(s, a)$ with a neural network parameterized by $\omega$. The network parameters $\omega$ can be updated with the TD algorithm. More specifically, the loss function is the weighted squared TD error averaged over the sampled mini-batch $\mathcal{B}_t$, i.e.,

$$ L(\omega) = \frac{1}{2|\mathcal{B}_t|} \sum_{e_i \in \mathcal{B}_t} w_i \delta_i^2, \tag{13} $$

where $|\mathcal{B}_t|$ represents the batch size and the importance-sampling weight $w_i$ is used to eliminate the bias introduced by prioritized sampling and ensure the same learning rate of all samples. According to [?], $w_i$ is given by

$$ w_i = \frac{(|\mathcal{B}_t| p_i)^{-\beta}}{\max_i w_i}, \tag{14} $$

which $\beta \in [0, 1]$ is a hyperparameter that controls the extent of the correction. Then, the gradient descent algorithm is leveraged to update the network parameters $\omega$ as

$$ \omega_{new} \leftarrow \omega_{now} - \alpha \nabla_\omega L(\omega_{now}), \tag{15} $$

where $\nabla_\omega L(\omega) = \frac{1}{|\mathcal{B}_t|} \sum_{e_i \in \mathcal{B}_t} w_i \delta_i \nabla_\omega Q(s_i, a_i; \omega)$ refers to the gradient of the loss function with respect to (w.r.t.) $\omega$, and $\alpha$ is the learning rate.

*3) Policy Network:* The policy network $\mu(s_t; \theta)$ aims to learn action policy by mapping the states to the specific actions. Since the action-value function $Q_\pi(s, a)$ can evaluate the score of the current action policy, the performance objective for $\mu(s_t; \theta)$ can be defined as [?]

$$ J(\theta) = \frac{1}{|\mathcal{B}_t|} \sum_{e_i \in \mathcal{B}_t} Q(s_i, \mu(s_i; \theta); \omega_{now}). \tag{16} $$

To learn the best policy, the parameters of the policy network can be optimized through the maximization of $J(\theta)$. Accordingly, the gradient ascend method is used to update $\theta$, i.e.,

$$ \theta_{new} \leftarrow \theta_{now} + \upsilon \nabla_\theta J(\theta_{now}), \tag{17} $$

where $\upsilon$ is the learning rate, and using chain rule yields $\nabla_\theta J(\theta) = \frac{1}{|\mathcal{B}_t|} \sum_{e_i \in \mathcal{B}_t} \nabla_\theta \mu(s_i; \theta) \nabla_a Q(s_i, a_i; \omega_{now})$.

*4) Target Evaluation/Policy Networks:* To further improve the stability, the soft update strategy is applied to update the parameters of the target networks, i.e., $\omega^-$ and $\theta^-$. More specifically, with the new parameters $\omega_{new}$ and $\theta_{new}$ given by (??) and (??), respectively, the parameters of the two target networks will be updated as

$$ \omega_{new}^- \leftarrow \tau \omega_{new} + (1 - \tau) \omega_{now}^-, \tag{18} $$

$$ \theta_{new}^- \leftarrow \tau \theta_{new} + (1 - \tau) \theta_{now}^-, \tag{19} $$

where the hyperparameter $\tau \ll 1$.

## IV. SIMULATIONS AND DISCUSSIONS

In this section, simulated results are presented for verification and discussion. For the system, the power is assumed equally allocated for XP-HARQ, i.e., $P_1 = \cdots = P_K$, and the average transmit signal-to-noise ratio (SNR) is defined as $P_1/\sigma^2 = \cdots = P_K/\sigma^2 \triangleq$ SNR. To deploy the DDPG, both the actor and critic networks consist of one input layer, three hidden layers, and one output layer. The number of the neurons in the three hidden layers are 100, 50, and 30 neurons, respectively. The three hidden layers of both networks use *ReLu* activation functions. The output layer of the actor network invokes *sigmoid* activation function to restrict the transmission rate within $R_k$, while the critical network does not leverage any activation function in the output layer. Both of the actor and critical networks capitalize on the adaptive moment estimation (Adam) optimizer to update the network parameters, and the learning rates are set to $\upsilon = \alpha = 0.001$. Furthermore, we assume that the number of epochs in the training state is 100, the number of time slots in each epoch is 6000, the size of the experienced reply buffer is $|\mathcal{M}| = 20000$,

the mini-batch size is 512. In addition, we assume that the weight of the soft update $\tau = 0.01$, the discount factor $\gamma = 0.9$, the extent of the correction $\beta = 0.5$, and the noise variance of the behavior policy $\vartheta^2 = 0.2$.

For illustration, the system parameters are set as $\sigma^2 = 1$, $\rho = 0.4$, and $\bar{R} = 10$ bps/Hz unless otherwise specified. Besides, we assume equal power allocation for XP-HARQ, i.e., $P_1 = \cdots = P_K$, and the average transmit signal-to-noise ratio (SNR) is defined as $P_1/\sigma^2 = \cdots = P_K/\sigma^2 \triangleq$ snr. To deploy the DDPG, both the actor and critic networks consist of one input layer, three hidden layers, and one output layer. The number of the neurons in the three hidden layers are 100, 50, and 30 neurons, respectively. In the meantime, the three hidden layers of both networks use "ReLu" activation functions. The output layer of the actor network invokes "sigmoid" activation function to restrict the transmission rate within $\bar{R}$, while the critical network does not leverage any activation function in the output layer. Both the actor and critical networks capitalize on the adaptive moment estimation (Adam) optimizer to update the network parameters, and the learning rates are set to $\mu = \epsilon_0 = 0.001$. Furthermore, we assume that the number of epochs in the training state is 100, the number of time slots in each epoch is 6000, the size of the prioritized replay buffer is $|\mathcal{M}| = 20000$, the mini-batch size is $|\mathcal{B}| = 512$.



Fig. 3. Impact of correlation coefficient on LTAT for different HARQ schemes.
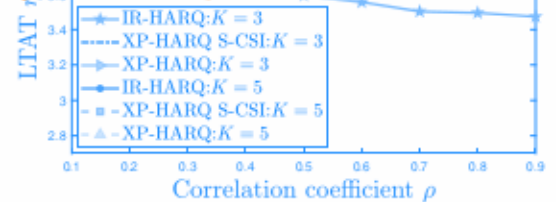


Fig. 4. Impact of correlation coefficient $\rho$.

## V. Conclusion

Due to the lack of simple analytical results of the performance metrics of XP-HARQ, we applied the DRL to properly select the incremental information rate for XP-HARQ over correlated fading channels, without recourse to the traditional optimization tools. More specifically, the maximization of the LTAT was formulated as a problem of MDP, which can be solved by using the algorithm of DDPG with prioritized experience replay. To demonstrate the efficacy of the proposed XP-HARQ scheme, its LTAT performance was compared to the conventional HARQ-IR and the XP-HARQ with only statistical CSI through simulations. It was found that IR-HARQ is more aggressive than XP-HARQ when determining the initial rate. In the meantime, it was also found that the time correlation has