

Introduction to Quantum Centralized-Critic and Multiple-Actor Networks

Joongheon Kim

School of Electrical Engineering, Korea University, Seoul, Republic of Korea

E-mail: joongheon@korea.ac.kr

Abstract—This paper introduces quantum centralized-critic and multiple-actor network architectures and organization.

I. INTRODUCTION

Quantum Centralized-Critic and Multiple-Actor Network.

We propose a novel quantum centralized critic and multiple actor network. Here, there is one *centralized critic* and multiple *actor* networks where the number of actor networks is commensurate with the number of agents. The agents make the sequential decision dispersively and train their multiple *actor* networks corresponding to the policy by evaluating the value of *centralized critic* network, which can be expressed as [1],

$$Q(o, a; \theta) = \beta_a \langle O_a \rangle_{o, \theta} = \beta_a \text{Tr}(U^{a\dagger}(o; \theta) M_a U^a(o; \theta)) \quad (1)$$

$$V(s; \phi) = \beta_c \langle O \rangle_{s, \phi} = \beta_c \text{Tr}(U^{c\dagger}(s; \phi) M_c U^c(s; \phi)) \quad (2)$$

where operators $\text{Tr}(\cdot)$, $U(\cdot)$, and $(\cdot)^\dagger$ represent trace operator, the unitary operation for qubit rotation, and the entanglement of multiple qubits and complex conjugate, respectively. When the quantum state is measured, the output (known as observable) exists between -1 and 1, *i.e.*, $\forall \langle O \rangle \in [-1, 1]$, we utilize hyper-parameters (β_a, β_c) for *actors* and *critic* networks to be well-trained. Note that M_a and M_c are Hermitian matrices. With Eqs. (1)–(2), and the hyper-parameters, the *actor-critic* networks approximate the value function.

Quantum Actor. At every time step t , the m -th quantum *actor* chooses the action with the most significant probability among the currently possible actions based on its state and observation information, which is represented as, $a_{m, \text{real}} = \arg \max D_a \pi_{\theta_m}(a_{\text{ideal}} | s_{\text{ideal}} + n_s, o_m) + n_a$, subject to $\pi_{\theta_m}(a_{\text{ideal}} | s_{\text{ideal}} + n_s, o_m) \triangleq \text{softmax}(Q(o, a; \theta_m))$ and $\text{softmax}(\mathbf{x}) \triangleq \left[\frac{e^{x_1}}{\sum_{i=1}^N e^{x_i}}, \dots, \frac{e^{x_N}}{\sum_{i=1}^N e^{x_i}} \right]$, where the $\text{softmax}(\cdot)$ is an activation function to normalize the inputs. By using it, we extract all actions' probabilities of the *actor* with the observable $\langle O_a \rangle_{o, \theta}$ in (1).

Quantum Centralized Critic. The CTDE has a *centralized critic* responsible for valuing the current state with a state-value function as follows, $V_\phi(s) = \langle O \rangle_{s, \phi} \simeq \mathbb{E}_{s_{\text{real}} \sim E, a_{\text{real}} \sim \pi_\theta} \left[\sum_{t'=t}^T \gamma^{t'-t} \cdot r(s_{\text{real}}, a_{\text{real}}, s'_{\text{real}}) \right]$, where s_t is the measured state at the current state at t . We also use the *critic* network's observable to evaluate the current state's value.

1) *Training and Inference:* MARL agents try to maximize the expected return. We utilize the congruent state-value function V_ϕ of the *centralized critic* network to derive the gradients from maximizing the common goal. With the parameters of

actor and *critic* networks, which correspond to θ and ϕ , we configure a multi-agent policy gradient (MAPG) based on the temporal difference *actor-critic* model by *Bellman optimality equation*, as follows,

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{s_{\text{real}}, o \sim E} \left[\sum_{t=1}^T \sum_{m=1}^M \delta_\phi^t \cdot \nabla_{\theta_m} \log \pi_{\theta_m}(a_m^t | s^t, o_m^t) \right], \quad (3)$$

and

$$\nabla_{\phi} \mathcal{L}(\phi) = \sum_{t=1}^T \nabla_{\phi} \|\delta_\phi^t\|^2, \quad (4)$$

subject to

$$\delta_\phi^t = r(s^t, a^t, s^{t+1}) + \gamma V_\phi(s^{t+1}) - V_\phi(s^t). \quad (5)$$

Among the three equations presented above, (3) is the objective function for *actor* networks, and the neural network parameters used in the equation are updated to be maximized by gradient ascent as, $\theta_m^{t+1} \approx \theta_m^t + \alpha_{\text{actor}} \times [\delta_\phi^t \cdot \nabla_{\theta} \log \pi_{\theta_m}(a_m^t | s^t, o_m^t)]$, where α_{actor} stands for a learning rate of *actor* networks. (4) corresponds to the *centralized critic* network's loss function which should be minimized by gradient descent as follows, $\phi^{t+1} \approx \phi^t + \alpha_{\text{critic}} \times [\delta_\phi^t \cdot \nabla_{\phi} V_\phi(s^t)]$, where α_{critic} is *centralized critic* network's learning rate. We describe how to obtain loss gradients with quantum and classical computing. Hereafter, we denote an *actor*-network and *critic* network as θ and ϕ for mathematical amenability. Their loss values are calculated with the temporal difference error of *centralized critic* δ_ϕ in (5), where the derivative of *actor/critic*'s i -th parameters is expressed as follows,

$$\frac{\partial J(\theta)}{\partial \theta_i} = \frac{\partial J(\theta)}{\partial \pi_\theta} \cdot \frac{\partial \pi_\theta}{\partial \langle O \rangle_{o, \theta}} \cdot \frac{\partial \langle O \rangle_{o, \theta}}{\partial \theta_i}, \quad (6)$$

$$\frac{\partial \mathcal{L}(\phi)}{\partial \phi_i} = \frac{\partial \mathcal{L}(\phi)}{\partial V_\phi} \cdot \frac{\partial V_\phi}{\partial \langle O \rangle_{s, \phi}} \cdot \frac{\partial \langle O \rangle_{s, \phi}}{\partial \phi_i}, \quad (7)$$

where the first and second derivatives of RHS can be calculated by classical computing. However, the latter derivative cannot be calculated because the quantum state is unknown before its measurement, *i.e.*, parameter-shift rule [1] is used.

REFERENCES

- [1] C. Park, W. J. Yun, J. P. Kim, T. K. Rodrigues, S. Park, S. Jung, and J. Kim, "Quantum multi-agent actor-critic networks for cooperative mobile access in multi-uav systems," *IEEE Internet of Things Journal*, pp. 1–1, 2023 (Early Access).