

# **Background: Reinforcement learning**

#### Goal

 Acquire an action-selection policy that maximizes a long-term reward by taking actions and observing the environment

• Q-value

State *s* Reward *r* 

Actor / Agent

Action

Expected value for action *a* in state *s*

- Q-learning algorithm
  - Q(s, a) is updated by taking action  $a_t$  and observing the next state  $s_{t+1}$  and reward  $r_t$

$$Q^{new}(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left\{ \begin{array}{c} r_t + \gamma \max_{a \in A} Q(s_{t+1}, a) - Q(s_t, a_t) \right\}$$

$$(1) \quad \text{Reward} \quad \text{Expected future value} \quad \text{Old value} \quad 2$$

### **Background: Deep Q-Network**

 Q-learning Estimated optimal future value  $Q(s_t, a_t) = r_t + \gamma \max_{a \in A} Q(s_{t+1}, a)$ Set of all the possible actions

- DQN: Q-learning + Deep neural network
  - Q-value is approximated with deep neural network  $\theta$
  - $-\theta$  is updated to minimize the loss

$$Q(s_t, a_t; \theta) = r_t + \gamma \max_{a \in A} Q(s_{t+1}, a; \theta)$$
$$L(\theta) = \left\{ r_t + \gamma \max_{a \in A} Q(s_{t+1}, a; \theta^-) - Q(s_t, a_t; \theta) \right\}^2$$
Loss function Fixed

– Game AI, robot control, communication control, …

 Various techniques are used for stability and convergence

# **Background: DQN techniques**

- Experience replay
  - Remove correlations in the observation sequence
  - Reuse experiences to increase sampling efficiency



- Distributed deep reinforcement learning
  - Generate more experiences using many machines
  - Ape-X DQN [1] (Distributed prioritized experience replay)

# Background: Ape-X

• Distributed prioritized experience replay [1]



### **1. Pull parameters** $\theta$ **2. Prediction**

Action:

 $a_t \leftarrow \varepsilon - \operatorname{greedy}(A)$ 

Next state & reward:

 $s_{t+1}, r_t \leftarrow \text{Environment}(a_t, s_t)$ 

LocalBuffer.ADD( $s_t, a_t, r_t, s_{t+1}$ )

[1] D. Horgan, et al., "Distributed Prioritized Experience Replay", ICLR'18.

 $\begin{array}{l} \text{if LocalBuffer.Size} \geq \text{BatchSize:} \\ \text{Priority:} \\ |\delta_t| \leftarrow |r_t + \gamma \max_{a \in A} Q(s_{t+1}, a) \\ -Q(s_t, a_t)| \end{array}$ 

#### 3. Push experience

 $\{a_t, s_t, r_t, s_{t+1}, \delta_t\} \times \text{BatchSize}$ 

a: Action, s: State, r: Reward,  $\delta$ : TD-error ,  $\gamma$ : Discount rate

5

### **Background: Ape-X**

• Distributed prioritized experience replay [1]



#### **1. Experience sampling**

 $\{a_t, s_t, r_t, s_{t+1}, \delta_t\} \times \text{BatchSize}$ 

#### 2. Batch training

$$Q(s_t, a_t; \boldsymbol{\theta}) \leftarrow r_t + \gamma \max_{a \in A} Q(s_{t+1}, a; \boldsymbol{\theta})$$

#### 3. Update priorities

#### **4.** Set parameters $\theta$

[1] D. Horgan, et al., "Distributed Prioritized Experience Replay", ICLR'18.

*a*: Action, *s*: State, *r*: Reward,  $\delta$ : TD-error,  $\gamma$ : Discount rate

6

# Background: Prioritized exp. replay

- Stochastic sampling [2]
  - Priority:  $p_t = |\delta_t| = |r_t + \gamma \max_{a \in A} Q(s_{t+1}, a) Q(s_t, a_t)|$
  - Sampling probability for experience *i*

 $P_{i} = \frac{p_{i}^{\alpha} - How \text{ much prioritization is used}}{\sum_{k} p_{k}^{\alpha}} \quad (p_{i} \neq 0)$ 

- Sum-tree structure
  - Add/Sample experience: O(log N)



# **Baseline: Our DRL implementation**

- Distributed deep reinforcement learning (DRL)
  - Actors, shared memory (Redis), and learner



3. Push parameters (13MB)

- 2. Prediction
- 3. Push experiences (43MB)

[3] Z. Wang, et al., "Dueling Network Architectures for Deep Reinforcement Learning", ICML'16.

### **Baseline: Our DRL implementation**

#### • Distributed deep reinforcement learning (DRL)



#### **Baseline: Execution time breakdown**

#### • Distributed deep reinforcement learning (DRL)



# **Baseline: Shared mem throughput**

- Distributed deep reinforcement learning (DRL)
- Actor push frequency
   Increase as # of actors
- Increase of # of actors
  - Network traffic increases
  - Gap between ideal and actual frequency increases



	Actors machine	Shared/Replay memory machine	Learner machine
OS	Ubuntu 20.04	Ubuntu 20.04	Ubuntu 20.04
CPU	Intel Xeon E5-2637 v3 @3.5GHz	Intel Xeon CPU E5-2637 v3 @3.50GHz	Intel Xeon CPU E5-1620 v2 @3.50GHz
Memory	128 GB	512 GB	128 GB
GPU	GeForce RTX 3080 ×1	—	GeForce GTX 1080Ti $\times 2$
CUDA/PyTorch	11.3 / 1.8.0+cu111	11.3 / 1.8.0+cu111	11.3 / 1.8.0+cu111
NIC	Intel Ethernet CNA XL710-QDA2	Intel Ethernet CNA XL710-QDA2	Intel Ethernet CNA XL710-QDA2
DPDK	_	20.11.0	_
Redis	_	5.0.5	
Network		Mellanox 40G Switch SX1012	

## Network optimizations on DRL

- Our approach
  - Accelerating shared memory & experience replay using DPDK
- DPDK (Data Plane Development Kit)
  - Network processing at application layer
  - Polling by dedicated CPU cores



# **Optimization 1: Shared memory**

In-network shared memory (Redis) by DPDK
 Low-latency shared memory by DPDK and F-stack



# **Optimization 1: Shared memory**

In-network shared memory (Redis) by DPDK



Redis access latency: 32.7% to 58.9% decrease

Replay memory pull throughput: 21.9% to 31.9% increase



# **Optimization 2: Experience replay**

• In-network experience sampling by DPDK

Moved from leaner machine to in-network switch



# **Optimization 2: Experience replay**

- In-network experience sampling by DPDK
  - Moved from leaner machine to in-network switch



# **Optimization 2: Experience replay**

• In-network experience sampling by DPDK

Moved from leaner machine to in-network switch



without F-Stack

with F-Stack

### Future work: Deployed in edge server

- Experience sampling in edge server
  - Actors located in edge, and learner located in cloud



# **Summary: Optimizations in DRL**

- Communication cost reduction by DPDK+F-stack
  - Optimization 1: In-network shared memory (Redis)
  - Optimization 2: In-network experience sampling



• Future work: Deployment in edge-cloud system

