



# HYPERDQN: A RANDOMIZED EXPLORATION METHOD FOR DEEP REINFORCEMENT LEARNING

#### Ziniu Li<sup>1</sup>, Yingru Li<sup>1†</sup>, Yushun Zhang<sup>1</sup>, Tong Zhang<sup>2†</sup>, and Zhi-Quan Luo<sup>1†</sup>

<sup>1</sup>Shenzhen Research Institute of Big Data, The Chinese University of Hong Kong, Shenzhen, China
<sup>2</sup>Hong Kong University of Science and Technology
{ziniuli, yingruli, yushunzhang}@link.cuhk.edu.cn,
tongzhang@ust.hk, luozq@cuhk.edu.cn

#### ICLR 2022





Problem: Exploration and Exploitation



# **Classic exploration methods**



1.ε-greedy : Select greedy action with probability of (1-ε)
 Select action randomly with probability of ε
 ε decreases with the increase of the agent step

2.Boltzmann distribution : Select action according to the Boltzmann distribution of Q-value (softmax)

3.Upper confidence bounds(UCB) : 
$$A_t^{UCB} \doteq arg \max_a \left[ Q_t + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

 $N_t(a)$ : The number of aciton a selected before time t

4. Thomson Sample : posterior sampling

(1)Set up Q-value distribution estimates for each action

(2)Sample randomly from each distribution to get Q-value, and select the action with the maximum Q-value

(3)Update distribution parameters based on reward

# **Benefits of posterior sampling**



• The upper regret bound is reduced

Regret $(T, \mathcal{M}) = \sum_{l=0}^{T/H-1} \mathbb{E}_{\mathcal{M}} \left[ V_0^*(s_{l0}) - \sum_{h=0}^{H} r_{lh} \right]$ 

• The randomness in posterior sampling could yield positive bias, which boosts optimistic behaviors

• the property of temporal extended exploration(deep exploration)

### Randomized least-squares value iteration (RLSVI)



Algorithm 2 RLSVI with greedy action **Input:** Features  $\Phi_0, ..., \Phi_{H-1}$ ;  $\sigma > 0, \lambda > 0$ 1: for l = 0, 1, ... do Compute  $\hat{\theta}_{l0}, ..., \hat{\theta}_{l,H-1}$  using Algorithm 1 2: Observe  $s_{l0}$ 3: for h = 0, ..., H - 1 do 4: Sample  $a_{lh} \in \operatorname{argmax}_{\alpha \in \mathcal{A}} \left( \Phi_h \tilde{\theta}_{lh} \right) (s_{lh}, \alpha)$ 5: Observe  $r_{lh}$  and  $s_{l,h+1}$ 6: end for 7: Observe  $r_{lH}$ 8: 9: **end for** 

Step1: Given initial value of  $\theta$  (l=0) Step2: Q-value comes from  $\Phi\theta$ ,and select action with greedy strategy

Step3: Update  $\theta$  according to Bayesian regression Step4: Repeat steps 2 and 3 for each episode Algorithm 1 Randomized Least-Squares Value Iteration Input: Data  $\Phi_0(s_{i0}, a_{i0}), r_{i0}, ..., \Phi_{H-1}(s_{iH-1}, a_{iH-1}), r_{iH}$ : i < L, Parameters  $\lambda > 0, \sigma > 0$ Output:  $\tilde{\theta}_{l0}, ..., \tilde{\theta}_{l,H-1}$ 1: for h = H - 1, ..., 1, 0 do 2: Generate regression problem  $A \in \mathbb{R}^{l \times K}, b \in \mathbb{R}^{l}$ :  $A \leftarrow \begin{bmatrix} \Phi_h(s_{0h}, a_{0h}) \\ \vdots \\ \Phi_h(s_{l-1,h}, a_{l-1,h}) \end{bmatrix}$  $b_i \leftarrow \begin{cases} r_{ih} + \max_{\alpha} (\Phi_{h+1} \tilde{\theta}_{l,h+1})(s_{i,h+1}, \alpha) & \text{if } h < H - 1 \\ r_{ih} + r_{i,h+1} & \text{if } h = H - 1 \end{cases}$ 

3: Bayesian linear regression for the value function

$$\overline{\theta}_{lh} \leftarrow \frac{1}{\sigma^2} \left( \frac{1}{\sigma^2} A^{\top} A + \lambda I \right)^{-1} A^{\top} b$$
$$\Sigma_{lh} \leftarrow \left( \frac{1}{\sigma^2} A^{\top} A + \lambda I \right)^{-1}$$

4: Sample  $\tilde{\theta}_{lh} \sim N(\bar{\theta}_{lh}, \Sigma_{lh})$  from Gaussian posterior 5: end for

### **Bayesian Linear Regression**



conditions: $y = \langle \theta^*, x \rangle + \omega^*, \ \omega^* \sim N(0, \sigma_{\omega}^2)$  prior distribution  $p_0(\theta^*) \sim N(\overline{\theta}_p, \sigma_p^2)$  Dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  According to Bayes formula:

$$p\big(\theta^*|\boldsymbol{\mathcal{D}}\big) = \frac{p\big(\boldsymbol{\mathcal{D}}|\theta^*\big)p_0\big(\theta^*\big)}{p(\boldsymbol{\mathcal{D}})} \ \Rightarrow \ p\big(\theta^*|\boldsymbol{\mathcal{D}}\big) \propto p\big(\boldsymbol{\mathcal{D}}|\theta^*\big)p_0\big(\theta^*\big) \ \Rightarrow \ p\big(\boldsymbol{\mathcal{D}}|\theta^*\big) = p\big(Y|\theta^*,X\big) \sim N\big(X\theta^*,\sigma_{\omega}^2\big)$$

Conjugate prior: For some likelihood functions, its prior and posterior have the same distribution Multi – dimensional gaussian distribution:

$$f(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} (X-\mu)^T \Sigma^{-1} (X-\mu)\right\}, \quad \Sigma: \text{ Covariance matrix}$$

So we can assume that the distribution:  $p\left(\theta^* | \boldsymbol{\mathcal{D}}\right) \sim N(\mu, \Sigma)$ 

For an ordinary multidimensional Gaussian distribution:

$$\begin{aligned} -\frac{1}{2}(X-\mu)^{T}\Sigma^{-1}(X-\mu) &= -\frac{1}{2}\left(X^{T}\Sigma^{-1}X-X^{T}\Sigma^{-1}\mu-\mu^{T}\Sigma^{-1}X+\mu^{T}\Sigma^{-1}\mu\right) \\ p(\mathcal{D}|\theta^{*}) &\propto p(\mathcal{D}|\theta^{*}) p_{0}(\theta^{*}) \propto \exp\left\{-\frac{1}{2\sigma_{\omega}^{2}}\left(Y^{T}-\theta^{*T}X^{T}\right)\left(Y-X\theta^{*}\right)\right\} \exp\left\{-\frac{1}{2\sigma_{p}^{2}}\left(\theta^{*}-\overline{\theta}_{p}\right)^{T}\left(\theta^{*}-\overline{\theta}_{p}\right)\right\} \\ &= \exp\left\{-\frac{1}{2\sigma_{\omega}^{2}}\left(Y^{T}Y-\theta^{*T}X^{T}Y-Y^{T}X\theta^{*}+\theta^{*T}X^{T}X\theta^{*}\right)-\frac{1}{2\sigma_{p}^{2}}\left(\theta^{*}-\overline{\theta}_{p}\right)^{T}\left(\theta^{*}-\overline{\theta}_{p}\right)\right\} \\ &= \exp\left\{-\frac{1}{2\sigma_{\omega}^{2}}\left(Y^{T}Y-\theta^{*T}X^{T}Y-Y^{T}X\theta^{*}+\theta^{*T}X^{T}X\theta^{*}\right)-\frac{1}{2}\left(\theta^{*}-\overline{\theta}_{p}\right)^{T}\left(\theta^{*}-\overline{\theta}_{p}\right)\right\} \\ &= \exp\left\{-\frac{1}{2\sigma_{\omega}^{2}}\left(Y^{T}Y-\theta^{*T}X^{T}Y-Y^{T}X\theta^{*}+\theta^{*T}X^{T}X\theta^{*}\right)-\frac{1}{2}\left(\theta^{*}-\overline{\theta}_{p}\right)^{T}\left(\theta^{*}-\overline{\theta}_{p}\right)\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\theta^{*T}\left(\sigma_{\omega}^{-2}X^{T}X+\sigma_{p}^{-2}I\right)\theta^{*}\right]+\frac{1}{2}\left[\theta^{*T}\left(\sigma_{\omega}^{-2}X^{T}Y+\sigma_{p}^{-2}\overline{\theta}_{p}\right)+\left(Y^{T}X\sigma_{\omega}^{-2}+\overline{\theta}_{p}^{T}\sigma_{p}^{-2}\right)\theta^{*}\right]+c\right\} \\ \text{Therefore} \qquad \Sigma = \sigma_{\omega}^{-2}X^{T}X+\sigma_{p}^{-2}I \quad \mu = \left(\sigma_{\omega}^{-2}X^{T}X+\sigma_{p}^{-2}I\right)\left(\sigma_{\omega}^{-2}X^{T}Y+\sigma_{p}^{-2}\overline{\theta}_{p}\right) \end{aligned}$$

# **Problem in extending to DRL**





Features Φ

Þ Q

> The computational complexity

fixed 
$$\phi$$
:  $\Phi_K = \Phi_{K-1} + \phi(x_K)\phi(x_K)^\top$  with  $\Phi_0 = \mathbf{0}$ ,  
changing  $\phi_K$ :  $\Phi_K := \sum_{\ell=1}^K \phi_K(x_\ell)\phi_K(x_\ell)^\top$ ,  $\Phi_{K-1} := \sum_{\ell=1}^{K-1} \phi_{K-1}(x_\ell)\phi_{K-1}(x_\ell)^\top$ ,  $\cdots$   
Covariance matrix

Hypermodel





**Theorem 1** Let  $p_z$  be the unit Gaussian distribution in  $\Re^K$ . For all  $\epsilon > 0$ ,  $\delta > 0$ , B > 0, and probability measures  $\mu$  over  $L_{\infty}(\mathcal{X}, B)$ , there exists a transport map H from  $p_z$  to  $\mu$ , a neural network  $f_{\theta} : \mathcal{X} \mapsto \Re$  with a linear output node and ReLU hidden nodes, and a linear hypermodel  $g_{\nu} : \mathcal{Z} \mapsto \Re^{N_{\theta}}$  such that

$$\|f_{g_{\nu}(z)} - f^*\|_{\infty} \le \epsilon,$$

with probability at least  $1 - \delta$ , for some  $\nu \in \Re^{N_{\nu}}$ , where  $f^* = H(z)$ .

A linear hypermodel can represent essentially any distribution over functions

HyperDQN





Loss function



$$L(\nu; \mathcal{D}) = \int_{z} p(z) \left[ \sum_{(x,y,\xi)\in\mathcal{D}} \left( y + \underbrace{\sigma_{\omega} z^{\mathsf{T}} \xi}_{(a)} - \underbrace{\left( g_{f_{\nu_{\text{prior}}}(z)}(x) + g_{f_{\nu}(z)}(x) \right)}_{(b)} \right)^{2} + \underbrace{\sigma_{\omega}^{2}}_{\sigma_{p}^{2}} \|f_{\nu}(z)\|^{2}}_{(c)} \right] (dz),$$

$$(4.1)$$

$$\min_{\nu,\theta_{\text{hidden}}} \int_{z} p(z) \left[ \sum_{(s,a,r,\xi,s')\in\mathcal{D}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{\mathsf{T}} \xi - Q_{\text{prediction}}(s,a,z) \right)^{2} + \frac{\sigma_{\omega}^{2}}{\sigma_{p}^{2}} \|f_{\nu}(z)\|^{2} \right] (dz),$$

$$(4.2)$$

where

$$Q_{\text{prediction}}(s, a, z) = Q_{\theta_{\text{prior}}, f_{\nu_{\text{prior}}}(z)}(s, a) + Q_{\theta_{\text{hidden}}, f_{\nu}(z)}(s, a),$$

$$Q_{\text{target}}(s', z) = r + \gamma \max_{a'} \left[ Q_{\theta_{\text{prior}}, f_{\nu_{\text{prior}}}(z)}(s', a') + Q_{\bar{\theta}_{\text{hidden}}, f_{\bar{\nu}}(z)}(s', a') \right].$$
(4.3)

 $\begin{array}{l} p(z): \text{Gaussian distribution} \\ \xi: \text{a random vector independently sampled from the unit hypersphere} & \begin{array}{c} \text{Generate } \mathbf{x} \sim N(0, I) \\ \\ \sigma_{\omega} z^{T} \xi: \text{is an artificial noise term exerted on the label y} \end{array} \\ \\ \begin{array}{c} \sigma_{\omega}^{2} \\ \sigma_{p}^{2} \end{array} \| f_{v}(z) \|^{2}: \text{a regularization term} \end{array} \right.$ 

#### Theorem



$$L(\nu; \mathcal{D}) = \int_{z} p(z) \left[ \sum_{\substack{(x,y,\xi) \in \mathcal{D}}} \left( y + \underbrace{\sigma_{\omega} z^{\mathsf{T}} \xi}_{(a)} - \underbrace{\left( g_{f_{\nu_{\text{prior}}}(z)}(x) + g_{f_{\nu}(z)}(x) \right)}_{(b)} \right)^{2} + \underbrace{\frac{\sigma_{\omega}^{2}}{\sigma_{p}^{2}} \left\| f_{\nu}(z) \right\|^{2}}_{(c)} \right] (\mathrm{d}z),$$

$$(4.1)$$

Assumption 1. Suppose the data generation follows  $y = x^{\top}\theta^{*} + \omega^{*}$ ,  $\omega^{*} \sim \mathcal{N}(0, \sigma_{\omega}^{2})$  and the prior distribution over  $\theta^{*}$  is  $\mathcal{N}(\bar{\theta}_{p}, \sigma_{p}^{2}I)$ . Furthermore, assume the base model is linear, i.e.,  $g_{f_{\nu}(z)}(x) = x^{\top}f_{\nu}(z)$  and  $g_{f_{\nu_{prior}}(z)}(x) = x^{\top}f_{\nu_{prior}}(z)$ . Moreover, assume the hypermodel is also linear, i.e.,  $f_{\nu}(z) = \nu_{w}^{\top}z + \nu_{b}$  and  $f_{\nu_{prior}}(z) = \nu_{w}^{prior}z + \nu_{b}^{prior}$ ;  $\nu = (\nu_{w}, \nu_{b})$  and  $\nu_{prior} = (\nu_{w}^{prior}, \nu_{b}^{prior})$ .

**Theorem 1** (Formal statement). Under Assumption 1, set  $\nu_w^{prior} = \sigma_p I$  and  $\nu_b^{prior} = \bar{\theta}_p$ . Let  $\nu^* = (\nu^*_w, \nu^*_b)$  be the optimal solution of (4.1) conditioned on specific realizations of  $\xi$ , then  $\theta := f_{\nu_{prior}}(z) + f_{\nu^*}(z) \sim \mathcal{N}(\nu_b^{prior} + \nu_b^*, (\nu_w^{prior} + \nu_w^*)^\top (\nu_w^{prior} + \nu_w^*))$  with  $\nu_b^{prior} + \nu_b^* = \mathbb{E}[\theta^* \mid \mathcal{D}], \quad (\nu_w^{prior} + \nu_w^*)^\top (\nu_w^{prior} + \nu_w^*) = \operatorname{Cov}[\theta^* \mid \mathcal{D}] + \operatorname{err}(\Xi),$ 

where

$$\operatorname{err}(\Xi) := \operatorname{Cov}[\theta^{\star} \mid \mathcal{D}] \left( \frac{1}{\sigma_{\omega}^{2}} \sum_{(x,\xi) \neq (x',\xi') \in \mathcal{D}} x\xi^{\top}\xi' x'^{\top} + \frac{1}{\sigma_{\omega}\sigma_{p}} \sum_{(x,\xi) \in \mathcal{D}} (x\xi^{\top} + \xi x^{\top}) \right) \operatorname{Cov}[\theta^{\star} \mid \mathcal{D}]$$

Furthermore, the error term satisfies  $\mathbb{E}_{\Xi}[err(\Xi)] = 0$ .

### **Noise choice**





z-dependent noise is indispensable

# Algorithm



#### Algorithm 2 HyperDQN

1: agent step $n \leftarrow 0$ , train step $\ell \leftarrow 0$ . 2: for episode $k = 0, 1, 2, \cdots$ do 3: generate a random vector $z \sim \mathcal{N}(0, I)$ . 5: for stage $t = 0, 1, 2, \cdots, T - 1$ do 6: observe state $s_i$ . 7: take the greedy action $a \leftarrow \operatorname{argmax}_a Q_{\theta}(s_t, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 then 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{hidden}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if $\frac{1}{ \widetilde{\mathcal{Z}} } \left(\sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left(Q_{target}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{prediction}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{s,a,r,\xi,s' \in \widetilde{\mathcal{D}}} \left(Q_{target}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{prediction}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{s,a,r,\xi,s' \in \widetilde{\mathcal{D}}} \left(Q_{target}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{prediction}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{s,a,r,\xi,s' \in \widetilde{\mathcal{D}}} \left(Q_{target}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{prediction}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ $	0			
2: for episode $k = 0, 1, 2, \cdots$ do 3: generate a random vector $z \sim \mathcal{N}(0, I)$ . 4: instantiate the Q-value function $Q_{\theta}(s, a, z)$ . 5: for stage $t = 0, 1, 2, \cdots, T - 1$ do 6: observe state $s_t$ . 7: take the greedy action $a \leftarrow \operatorname{argmax}_a Q_{\theta}(s_t, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n, \tan$ frequency $M$ ) == 0 then 13: sample $n$ mandom vectors $z'_t: \tilde{Z} = \{z'_t\}_{t=1}^N$ . 14: sample $N$ random vectors $z'_t: \tilde{Z} = \{z'_t\}_{t=1}^N$ . 15: optimize $\nu$ and $\theta_{hidden}$ using the empirical loss function (A.6) with $\tilde{\mathcal{D}}$ and $\tilde{Z}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for 22: end for 23: end for 24: $\sum_{x \in \tilde{Z}} \frac{ \mathcal{D} }{ \tilde{Z} } \sum_{x \in \tilde{Z}} \frac{ \mathcal{D} }{ \tilde{D} } \sum_{(s,a,r,\xi,s') \in \tilde{\mathcal{D}}} \left( Q_{target}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{prediction}(s, a, z) \right)$	1: 8	agent step $n \leftarrow 0$ , train step $\ell \leftarrow 0$ .		
3: generate a random vector $z \sim \mathcal{N}(0, I)$ . instantiate the $Q$ -value function $Q_{\theta}(s, a, z)$ . 5: for stage $t = 0, 1, 2, \dots, T - 1$ do beserve state $s_t$ . 7: take the greedy action $a \leftarrow \operatorname{argmax}_a Q_{\theta}(s_t, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 then 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z_i^t$ : $\widetilde{\mathcal{Z}} = \{z_i^t\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{hidden}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end for 21: end for 22: end for 22: end for 23: $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{target}(s', z) + \sigma_{\omega} z^{\top} \xi - Q_{prediction}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \frac{\xi}{ \widetilde{\mathcal{L}} ^2} + \frac{1}{ \widetilde{\mathcal{D}} } \frac{\xi}{ \widetilde{\mathcal{L}} ^2} + \frac{1}{ \widetilde{\mathcal{D}} } \frac{\xi}{ \widetilde{\mathcal{L}} ^2} + \frac{\xi}{ \widetilde{\mathcal{D}} ^2} + \frac{\xi}{ \widetilde{\mathcal{L}} ^2} + \frac{\xi}{ \mathcal{$	2: 1	for episode $k = 0, 1, 2, \cdots$ do		
4: Instantiate the $Q$ -value function $Q_{\theta}(s, a, z)$ . 5: for stage $t = 0, 1, 2, \dots, T - 1$ do 6: observe state $s_t$ . 7: take the greedy action $a \leftarrow \operatorname{argmax}_a Q_{\theta}(s_t, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n, \text{train frequency } M ) == 0$ then 13: sample $a$ mini-bach $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i$ : $\widetilde{\mathcal{Z}} = \{z_i^1\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for 22: end for 22: end for	3:	generate a random vector $z \sim \mathcal{N}(0, I)$ .	⊳ Sampling	
5: <b>for</b> stage $t = 0, 1, 2, \dots, T - 1$ <b>do</b> 6: observe state $s_t$ . 7: take the greedy action $a \leftarrow \arg argmax_a Q_{\theta}(s_t, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: <b>if</b> mod (agent step $n$ , train frequency $M$ ) == 0 <b>then</b> 13: sample a mini-bach $\tilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i$ : $\tilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{hidden}$ using the empirical loss function (A.6) with $\tilde{\mathcal{D}}$ and $\tilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: <b>end if</b> 18: <b>if</b> mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: <b>end if</b> 21: <b>end for</b> 22: <b>end for</b> 22: <b>end for</b> 22: <b>end for</b> 23: <b>end for</b> 24: $ \tilde{\mathcal{Z}} \left(\sum_{z \in \tilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \tilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \tilde{\mathcal{D}}} \left(Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z)\right)$	4:	instantiate the Q-value function $Q_{\theta}(s, a, z)$ .		
6: observe state $s_t$ . 7: take the greedy action $a \leftarrow \operatorname{argmax}_a Q_{\theta}(s_t, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 <b>then</b> $\triangleright$ Update 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: <b>end if</b> 18: <b>if</b> mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: <b>end if</b> 21: <b>end for</b> 22: <b>end for</b> 22: <b>end for</b> 22: <b>end for</b> 23: $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) + \frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \mathcal{$	5:	for stage $t = 0, 1, 2, \cdots, T - 1$ do	▷ Interaction	
7: take the greedy action $a \leftarrow \operatorname{argma}_{a} Q_{\theta}(s_{t}, a, z)$ . 8: receive the next state $s_{t+1}$ and reward $r(s_{t}, a_{t})$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_{t}, a_{t}, r, \xi, s_{t+1})$ into the replay biffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 <b>then</b> $\triangleright$ Update 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_{i}: \widetilde{\mathcal{Z}} = \{z'_{i}\}_{i=1}^{N}$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: end if 21: end for 22: end for 22: end for 23: end for	6:	observe state $s_t$ .		
8: receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ . 9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 <b>then</b> $\triangleright$ Update 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i$ : $\widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: end if 21: end for 22: end for 22: end for 22: end for 23: $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) \right)$	7:	take the greedy action $a \leftarrow \operatorname{argmax}_a Q_{\theta}(s_t, a, z)$ .		
9: sample $\xi$ uniformly from unit hypersphere. 10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 <b>then</b> $\triangleright$ Update 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: end if 21: end for 22: end for $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) \right)$	8:	receive the next state $s_{t+1}$ and reward $r(s_t, a_t)$ .		
10: store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ . 11: agent step $n \leftarrow n+1$ . 12: <b>if</b> mod (agent step $n$ , train frequency $M$ ) == 0 <b>then</b> $\triangleright$ Update 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z_i^t: \widetilde{\mathcal{Z}} = \{z_i^t\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: <b>end if</b> 18: <b>if</b> mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: <b>end if</b> 21: <b>end for</b> 22: <b>end for</b> $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) \right)$	9:	sample $\xi$ uniformly from unit hypersphere.		
11: agent step $n \leftarrow n + 1$ . 12: if mod (agent step $n$ , train frequency $M$ ) == 0 then $\triangleright$ Update 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for $\frac{1}{ \widetilde{\mathcal{Z}} } \left(\sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left(Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z)\right)$	10:	store $(s_t, a_t, r, \xi, s_{t+1})$ into the replay buffer $\mathcal{D}$ .		
12: <b>if</b> mod (agent step <i>n</i> , train frequency <i>M</i> ) == 0 <b>then</b> 13: sample a mini-batch $\widetilde{\mathcal{D}}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample <i>N</i> random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: <b>end if</b> 18: <b>if</b> mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: <b>end if</b> 21: <b>end for</b> 22: <b>end for</b> 22: <b>end for</b> 22: <b>end for</b> 1 $ \widetilde{\mathcal{Z}}  \left(\sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left(Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z)\right)$	11:	agent step $n \leftarrow n+1$ .		
13: sample a mini-batch $\mathcal{D}$ of $(s, a, r, \xi, s')$ from the replay buffer $\mathcal{D}$ . 14: sample $N$ random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for 22: end for $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s', z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s, a, z) \right)$	12:	if mod (agent step n, train frequency $M$ ) == 0 then	⊳ Update	
14: sample N random vectors $z'_i: \widetilde{Z} = \{z'_i\}_{i=1}^N$ . 15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{D}$ and $\widetilde{Z}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for 22: end for $\frac{1}{ \widetilde{Z} } \left( \sum_{z \in \widetilde{Z}} \frac{ \mathcal{D} }{ \widetilde{D} } \sum_{(s,a,r,\xi,s') \in \widetilde{D}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right)$	13:	sample a mini-batch $\mathcal{D}$ of $(s, a, r, \xi, s')$ from the replay by	Iffer $\mathcal{D}$ .	
15: optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A.6) with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ . 16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right)$	14:	sample N random vectors $z'_i: \widetilde{\mathcal{Z}} = \{z'_i\}_{i=1}^N$ .		
16: train step $\ell \leftarrow \ell + 1$ . 17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for $\frac{1}{ \widetilde{Z} } \left( \sum_{z \in \widetilde{Z}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	15:	optimize $\nu$ and $\theta_{\text{hidden}}$ using the empirical loss function (A	$\widetilde{\mathcal{D}}$ , with $\widetilde{\mathcal{D}}$ and $\widetilde{\mathcal{Z}}$ .	
17: end if 18: if mod (train step $\ell$ , target update frequency $G$ ) == 0 then 19: update the target network. 20: end if 21: end for 22: end for $\frac{1}{ \widetilde{Z} } \left( \sum_{z \in \widetilde{Z}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	16:	train step $\ell \leftarrow \ell + 1$ .		
18: <b>if</b> mod (train step $\ell$ , target update frequency $G$ ) == 0 <b>then</b> 19: update the target network. 20: <b>end if</b> 21: <b>end for</b> 22: <b>end for</b> $\frac{1}{ \widetilde{Z} } \left( \sum_{z \in \widetilde{Z}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	17:	end if		
19: update the target network. 20: end if 21: end for 22: end for $\frac{1}{ \widetilde{Z} } \left( \sum_{z \in \widetilde{Z}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	18:	if mod (train step $\ell$ , target update frequency $G$ ) == 0 then		
20: end if 21: end for 22: end for $\frac{1}{ \widetilde{Z} } \left( \sum_{z \in \widetilde{Z}} \frac{ \mathcal{D} }{ \widetilde{D} } \sum_{(s,a,r,\xi,s') \in \widetilde{D}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	19:	update the target network.		
21: end for 22: end for $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	20:	end if		
22: end for $\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	21:	end for		
$\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{T} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$	22: (	end for		
$\frac{1}{ \widetilde{\mathcal{Z}} } \left( \sum_{z \in \widetilde{\mathcal{Z}}} \frac{ \mathcal{D} }{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left( Q_{\text{target}}(s',z) + \sigma_{\omega} z^{\top} \xi - Q_{\text{prediction}}(s,a,z) \right) \right)$		/	∲	
$\overline{ \widetilde{\mathcal{Z}} } \left[ \sum_{z \in \widetilde{\mathcal{Z}}} \frac{1}{ \widetilde{\mathcal{D}} } \sum_{(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}} \left[ Q_{\text{target}}(s',z) + \sigma_{\omega} z + \xi - Q_{\text{prediction}}(s,a,z) \right] \right]$		$1 \int \sum  \mathcal{D} $	$\sum \left( \begin{array}{c} c \\ c \end{array} \right) \left( \begin{array}{c} c \end{array} \right) \left( \begin{array}{c} c \\ c \end{array} \right) \left( \begin{array}{c} c \\ c \end{array} \right) \left( \begin{array}{c} c \\ c \end{array} \right) \left( \begin{array}{c} c \end{array} \right) \left( \begin{array}{c} c \\ c \end{array} \right) \left( \begin{array}{c} c \end{array} \right) $	
$\begin{array}{c c}  \mathcal{L}  & \sum_{z \in \widetilde{\mathcal{Z}}}  \mathcal{L}  & (s, a, r, \xi, s') \in \widetilde{\mathcal{D}} \end{array} \\ \end{array} $		$\overline{ \widetilde{\mathfrak{z}} } \int \sum \overline{ \widetilde{\mathfrak{D}} }$	$\sum \qquad (Q_{\text{target}}(s',z) + \sigma_{\omega}z'\xi -$	- $Q_{\text{prediction}}(s, a, z)$
		$ \mathcal{Z}  = \int_{z \in \widetilde{\mathcal{Z}}}  \mathcal{D} $ (s	$(s,a,r,\xi,s') \in \widetilde{\mathcal{D}}$	/

### Experiment



#### Atari

Table 2: Comparison of algorithms on Atari in terms of the median over 49 games' maximum human-normalized scores. Note that the performance of DQN is based on 200M training frames while other methods are based on 20M training frames.

DQN (200M)	OPIQ	OB2I	BootDQN	NoisyNet	HyperDQN
93%	37%	50%	82%	91%	110%



#### Experiment



### SuperMarioBros

Table 3: Comparison of algorithms on SuperMarioBros in terms of the raw scores by the best policies with 20M training frames.

	DQN	OPIQ	OB2I	BootDQN	NoisyNet	HyperDQN
SuperMarioBros-1-1	1,070	7,650	4,457	7,009	12,439	7,924
SuperMarioBros-1-2	2,883	5,515	4,695	5,665	6,347	8,267
SuperMarioBros-1-3	667	2,053	1,583	1,609	1,587	<b>6,047</b>
SuperMarioBros-2-1	10,800	21,654	14,226	${\bf 26, 415}$	14,017	23,047
SuperMarioBros-2-2	813	1,630	1,588	1,092	1,808	1,984
SuperMarioBros-2-3	3,373	4,718	4,402	5,108	$\mathbf{6, 490}$	5,980
SuperMarioBros-3-1	2,560	3,700	3,251	3,862	11,310	${f 48, 385}$
SuperMarioBros-3-2	11,633	20,872	26,508	20,955	33,489	41, 140
SuperMarioBros-3-3	1,007	2,440	3,009	2,650	${\bf 5,886}$	5,568

#### Deep Sea





# Thanks