

---

# Distributional Reinforcement Learning via Sinkhorn Iterations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Distributional reinforcement learning (RL) is a class of state-of-the-art algorithms  
2        that estimate the whole distribution of the total return rather than only its expect-  
3        ation. The representation manner of each return distribution and the choice of  
4        distribution divergence are pivotal for the empirical success of distributional RL.  
5        In this paper, we propose a new class of *Sinkhorn distributional RL (Sinkhorn-*  
6        *DRL)* algorithm that learns a finite set of statistics, i.e., deterministic samples,  
7        from each return distribution and then leverages Sinkhorn iterations to evaluate  
8        the Sinkhorn distance between the current and target Bellman distributions. Re-  
9        markably, Sinkhorn divergence interpolates between the Wasserstein distance and  
10       Maximum Mean Discrepancy (MMD). This allows our proposed SinkhornDRL  
11       algorithm to find a sweet spot leveraging the geometry of optimal transport based  
12       distance and the unbiased gradient estimates of MMD. Finally, experiments on  
13       the suit of 55 Atari games reveal the competitive performance of SinkhornDRL  
14       algorithm as opposed to existing state-of-the-art algorithms.

## 15    1 Introduction

16    Classical reinforcement learning (RL) algorithms are normally based on the expectation of discounted  
17    cumulative rewards that an agent observes while interacting with the environment. Recently, a new  
18    class of RL algorithms called *distributional RL* estimates the full distribution of total returns and has  
19    exhibited the state-of-the-art performance in a wide range of environments [2, 8, 7, 24, 26, 17].

20    From the literature of distributional RL, it is easily recognized that algorithms based on either  
21    Wasserstein distance or MMD have gained great attention due to their superior performance. As such,  
22    their mutual connection from the perspective of mathematical properties intrigues us to explore further  
23    in order to design new algorithms. Particularly, Wasserstein distance, long known to be a powerful  
24    tool to compare probability distributions with non-overlapping supports, has recently emerged as an  
25    appealing contender in various machine learning applications. It is known that Wasserstein distance  
26    was long disregarded because of its computational burden in its original form to solve an expensive  
27    network flow problem. However, recent works [21, 14] have shown that this cost can be largely  
28    mitigated by settling for cheaper approximations through strongly convex regularizers. The benefit of  
29    this regularization has opened the path to wider applications of the Wasserstein distance in relevant  
30    learning problems, including the design of distributional RL algorithms.

31    The Sinkhorn divergence [21] introduces the entropic regularization on the Wasserstein distance,  
32    allowing it tractable for the evaluation especially in high-dimensions. It has been successfully applied  
33    in numerous crucial machine learning developments, including the Sinkhorn-GAN [14] and Sinkhorn-  
34    based adversarial training [23]. More importantly, it has been shown that Sinkhorn divergence  
35    interpolates Wasserstein distance and MMD, and their equivalence form can be well established in the  
36    limit cases [11, 18, 17]. However, a Sinkhorn-based distributional RL algorithm has not yet been

37 formally proposed and its connection with algorithms based on Wasserstein distance and MMD is  
 38 also less studied. Therefore, a natural question is *can we design a new class of distributional RL*  
 39 *algorithms via Sinkhorn divergence, thus bridging the gap between existing two main branches of*  
 40 *distributional RL algorithms?* Moreover, the dominant quantile-based algorithms, e.g., QR-DQN [8],  
 41 aimed at approximating Wasserstein distance, suffers from the non-crossing issue in the quantile  
 42 estimation [26], while sample-based Sinkhorn algorithm can naturally circumvent this problem.

43 In this paper, we propose a novel distributional RL algorithm based on *Sinkhorn divergence*. Firstly,  
 44 we point out the key roles of distribution divergence and representation of value distribution in the  
 45 design of distributional RL. After a detailed introduction of our proposed SinkhornDRL algorithm,  
 46 we theoretically analyze its convergence guarantee and moment matching behavior of distributional  
 47 Bellman operators under Sinkhorn divergence. Thus, a regularized MMD equivalence form of  
 48 Sinkhorn divergence is derived, interpreting the empirical success of our algorithms in real applications.  
 49 Finally, we compare the performance of our SinkhornRL algorithm with typical baselines on 55 Atari  
 50 games, verifying the competitive performance of our proposal. Our approach inspires researchers  
 51 to find a trade-off that simultaneously leverages the geometry of the Wasserstein distance and the  
 52 favorable unbiased gradient estimate property of MMD while designing new distributional RL  
 53 algorithms in the future.

## 54 2 Preliminary Knowledge

### 55 2.1 Distributional Reinforcement Learning

56 In the classical RL setting, an agent interacts with an environment via a Markov decision pro-  
 57 cess (MDP), a 5-tuple  $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$ , where  $\mathcal{S}$  and  $\mathcal{A}$  are the state and action spaces, respectively.  $P$   
 58 is the environment transition dynamics,  $R$  is the reward function and  $\gamma \in (0, 1)$  is the discount factor.

59 **From Value function to Value distribution.** Given a policy  $\pi$ , the discounted sum of future  
 60 rewards is a random variable  $Z^\pi(s, a) = \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)$ , where  $s_0 = s$ ,  $a_0 = a$ ,  $s_{t+1} \sim$   
 61  $P(\cdot | s_t, a_t)$ , and  $a_t \sim \pi(\cdot | s_t)$ . In the control setting, expectation-based RL is based on the action-  
 62 value function  $Q^\pi(s, a)$ , which is the expectation of  $Z^\pi(s, a)$ , i.e.,  $Q^\pi(s, a) = \mathbb{E}[Z^\pi(s, a)]$ . By  
 63 contrast, distributional RL focuses on the action-value distribution, the full distribution of  $Z^\pi(s, a)$ ,  
 64 and the incorporation of additional distributional knowledge intuitively interprets its empirical success.

65 **Distributional Bellman operators.** For the policy evaluation in expectation-based RL, the action-  
 66 value function is updated via the Bellman operator  $\mathcal{T}^\pi Q(s, a) = \mathbb{E}[R(s, a)] + \gamma \mathbb{E}_{s' \sim p, \pi} [Q(s', a')]$ .  
 67 In distributional RL, the action-value distribution of  $Z^\pi(s, a)$  is updated via the distributional Bellman  
 68 operator  $\mathfrak{T}^\pi$

$$\mathfrak{T}^\pi Z(s, a) = R(s, a) + \gamma Z(s', a'), \quad (1)$$

69 where  $s' \sim P(\cdot | s, a)$  and  $a' \sim \pi(\cdot | s')$ . The equality in Eq. 1 implies that random variables of  
 70 both sides are equal in distribution. The distributional Bellman operator  $\mathfrak{T}^\pi$  is contractive under  
 71 certain distribution divergence metrics, but the distributional Bellman optimality operator  $\mathfrak{T}$  can only  
 72 converge to a set of optimal non-stationary value distributions in a weak sense [9].

### 73 2.2 Divergences between Measures

74 **Optimal Transport (OT) and Wasserstein Distance** The optimal transport (OT) metric between  
 75 two probability measures  $(\mu, \nu)$  supported on two metric spaces is defined as the solution of the linear  
 76 program:

$$\min_{\Pi \in \Pi(\mu, \nu)} \int c(x, y) d\Pi(x, y), \quad (2)$$

77 where  $c$  is the cost function and  $\Pi$  is the joint distribution with marginals  $(\mu, \nu)$ . Wasserstein distance  
 78 (a.k.a. earth mover distance) is a special case of optimal transport with the Euclidean norm as the  
 79 cost function. In particular, given two scalar random variables  $X$  and  $Y$ ,  $p$ -Wasserstein metric  $W_p$   
 80 between the distributions of  $X$  and  $Y$  can be simplified as

$$W_p(X, Y) = \left( \int_0^1 |F_X^{-1}(\omega) - F_Y^{-1}(\omega)|^p d\omega \right)^{1/p}, \quad (3)$$



81 where  $F^{-1}$  is the inverse cumulative distribution function of a random variable. The desirable  
 82 geometric property of Wasserstein distance allows it to recover full support of measures, but it suffers  
 83 from the curse of dimension [13, 1].

84 **Maximum Mean Discrepancy** The squared Maximum Mean Discrepancy (MMD)  $\text{MMD}_k^2$  with  
 85 the kernel  $k$  is formulated as

$$\text{MMD}_k^2 = \mathbb{E}[k(X, X')] + \mathbb{E}[k(Y, Y')] - 2\mathbb{E}[k(X, Y)], \quad (4)$$

86 where  $k(\cdot, \cdot)$  is a continuous kernel on  $\mathcal{X}$ .  $X'$  (resp.  $Y'$ ) is a random variable independent of  $X$   
 87 (resp.  $Y$ ). If  $k$  is a trivial kernel, MMD degenerates to the energy distance. Mathematically, the “flat”  
 88 geometry that MMD induces on the space of probability measures does not faithfully lift the ground  
 89 distance [11], but MMD is cheaper to compute than OT and has a smaller sample complexity, i.e.,  
 90 approximating the distance with samples of measures [13]. We provide the detailed introduction of  
 91 more distribution divergences in Appendix A.

### 92 3 Roles of Distribution Divergence and Representation in distributional RL

#### 93 3.1 Distributional RL: From Neural Q-Fitted Iteration to Neural Z-Fitted Iteration

94 **Neural Q-Fitted Iteration.** It is known that Deep Q Learning [16] can be simplified into *Neural*  
 95 *Q-Fitted Iteration* [10] under tricks of experience replay and the target network  $Q_{\theta^*}$ , where we update  
 96 parameterized  $Q_{\theta}(s, a)$  in each iteration  $k$ :

$$Q_{\theta}^{k+1} = \underset{Q_{\theta}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n [y_i - Q_{\theta}^k(s_i, a_i)]^2, \quad (5)$$

97 where the target  $y_i = r(s_i, a_i) + \gamma \max_{a \in \mathcal{A}} Q_{\theta^*}^k(s'_i, a)$  is fixed within every  $T_{\text{target}}$  steps to update  
 98 target network  $Q_{\theta^*}$  by letting  $\theta^* = \theta$  and the experience buffer induces independent samples  
 99  $\{(s_i, a_i, r_i, s'_i)\}_{i \in [n]}$ . In an ideal case that neglects the non-convexity and TD approximation errors,  
 100 we have  $Q_{\theta}^{k+1} = \mathcal{T}Q_{\theta}^k$ , which is exactly equivalent to updating under Bellman optimality operator.

101 **Neural Z-Fitted Iteration.** Analogous to neural Q-fitted iteration, we can also simplify value-based  
 102 distributional RL methods based on a parameterized  $Z_{\theta}$  into a *Neural Z-fitted Iteration* as

$$Z_{\theta}^{k+1} = \underset{Z_{\theta}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n d_p(Y_i, Z_{\theta}^k(s_i, a_i)), \quad (6)$$

103 where the target  $Y_i = R(s_i, a_i) + \gamma Z_{\theta^*}^k(s'_i, \pi_Z(s'))$  with  $\pi_Z(s') = \operatorname{argmax}_{a'} \mathbb{E}[Z_{\theta^*}^k(s', a')]$  is  
 104 fixed within every  $T_{\text{target}}$  steps to update target network  $Z_{\theta^*}$ , and  $d_p$  is a divergence metric between  
 105 two distributions.

#### 106 3.2 Key Roles of $d_p$ and $Z_{\theta}$

107 Within the Neural Z-fitted Iteration framework proposed in Eq. 6, we observe that the choice of  
 108 representation manner on  $Z_{\theta}$  and the metric  $d_p$  are pivotal for the distributional RL algorithms. For  
 109 instance, QR-DQN [8] approximates Wasserstein distance  $W_p$ , which leverages quantiles to represent

Algorithm	$d_p$ Distribution Divergence	Representation $Z_{\theta}$	Convergence Rate of $\mathfrak{T}^{\pi}$	Sample Complexity of $d_p$
C51 [2]	Cramér distance	Histogram	$\sqrt{\gamma}$	$\searrow$
QR-DQN [8]	Wasserstein distance	Quantiles	$\gamma$	$\mathcal{O}(n^{-\frac{1}{d}})$
MMDDRL [17]	MMD	Samples	$\gamma^{\alpha/2}$ with kernel $k_{\alpha}$	$\mathcal{O}(1/n)$
SinkhornDRL (ours)	Sinkhorn divergence	Samples	$\gamma (\varepsilon \rightarrow 0)$ $\gamma^{\alpha/2} (\varepsilon \rightarrow \infty)$	$\mathcal{O}(n^{\frac{\varepsilon}{\varepsilon + d/2} \sqrt{\pi}}) (\varepsilon \rightarrow 0)$ $\mathcal{O}(n^{-\frac{1}{2}}) (\varepsilon \rightarrow \infty)$

Table 1: Comparison between typical distributional RL algorithms under different distribution divergences and representation of  $Z_{\theta}$ .  $k_{\alpha} = -\|x - y\|^{\alpha}$  in MMDDRL,  $d$  is the sample dimension and  $\kappa = 2\beta d + \|c\|_{\infty}$ , where the cost function  $c$  is  $\beta$ -Lipschitz [13]. Sample complexity of MMD can be improved to  $\mathcal{O}(1/n)$  using kernel herding technique [5].

110 the distribution of  $Z_\theta$ . C51 [2] represents  $Z_\theta$  via a categorical distribution under the convergence of  
 111 Cramér distance [3, 19], while MMD distributional RL (MMDDL) [17] learns samples to represent  
 112 the distribution of  $Z_\theta$  based on MMD. We compare characteristics of these distribution divergence,  
 113 including the convergence rate and sample complexity, in Table 1. Theoretical results regarding  
 114 Sinkhorn divergence is based on [13] and the detailed convergence proof of other distances is also  
 115 provided in Appendix A. In summary, we argue that  $d_p$  and  $Z_\theta$  are two crucial factors in distributional  
 116 RL design, based on which we introduce our Sinkhorn distributional RL.

## 117 4 Sinkhorn Distributional RL (SinkhornDRL)

118 In this section, we firstly introduce Sinkhorn divergence and apply it in distributional RL. Next, we  
 119 conduct a theoretical analysis about the convergence speed and a new moment matching manner of  
 120 our algorithm under the Sinkhorn divergence. Finally, a practical Sinkhorn iteration algorithm is  
 121 introduced to evaluate the Sinkhorn divergence.

### 122 4.1 Sinkhorn Divergence and Genetic Algorithm

123 We design Sinkhorn distributional RL algorithm via Sinkhorn divergence. Sinkhorn divergence [21] is  
 124 a tractable loss to approximate the optimal transport problem by leveraging an entropic regularization  
 125 to turn the original Wasserstein distance into a differentiable and more robust quantity. The resulting  
 126 loss can be computed using Sinkhorn fixed point iterations, which is naturally suitable for modern deep  
 127 learning frameworks. In particular, the entropic smoothing generates a family of losses interpolating  
 128 between Wasserstein distance and Maximum Mean Discrepancy (MMD). As such, it allows us to find  
 129 a sweet trade-off that simultaneously leverages the geometry of Wasserstein distance on the one hand,  
 130 and the favorable high-dimensional sample complexity and unbiased gradient estimates of MMD. We  
 131 introduce the entropic regularized Wasserstein distance  $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$  as

$$\min_{\Pi \in \Pi(\mu, \nu)} \int c(x, y) d\Pi(x, y) + \varepsilon \text{KL}(\Pi | \mu \otimes \nu), \quad (7)$$

132 where  $\text{KL}(\Pi | \mu \otimes \nu) = \int \log \left( \frac{\Pi(x, y)}{d\mu(x) d\nu(y)} \right) d\Pi(x, y)$  is a strongly convex regularization. The impact  
 133 of this entropy regularization is similar to  $\ell_2$  ridge regularization in linear regression. Next, the  
 134 Sinkhorn loss [11, 14] between two measures  $\mu$  and  $\nu$  is defined as

$$\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu) = 2\mathcal{W}_{c,\varepsilon}(\mu, \nu) - \mathcal{W}_{c,\varepsilon}(\mu, \mu) - \mathcal{W}_{c,\varepsilon}(\nu, \nu). \quad (8)$$

135 As demonstrated by [11], the Sinkhorn divergence  $\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu)$  is convex, smooth and positive definite  
 136 that metrizes the convergence in law. In statistical physics,  $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$  can be re-factored as a  
 137 projection problem:

$$\mathcal{W}_{c,\varepsilon}(\mu, \nu) := \min_{\Pi \in \Pi(\mu, \nu)} \text{KL}(\Pi | \mathcal{K}), \quad (9)$$

138 where  $\mathcal{K}$  is the Gibbs distribution with the density function satisfies  $d\mathcal{K}(x, y) = e^{-\frac{c(x, y)}{\varepsilon}} d\mu(x) d\nu(y)$ .  
 139 This problem is often referred to as the “static Schrödinger problem” [15, 20] as it was initially  
 140 considered in statistical physics.

141 **Distributional RL with Sinkhorn Divergence and Particle Representation.** The key of apply-  
 142 ing Sinkhorn divergence in distributional RL is to simply leverage the Sinkhorn loss  $\overline{\mathcal{W}}_{c,\varepsilon}$  to mea-  
 143 sure the distance between the current action-value distribution  $Z_\theta(s, a)$  and the target distribution  
 144  $\mathfrak{T}^\pi Z_\theta(s, a)$ , yielding  $\overline{\mathcal{W}}_{c,\varepsilon}(Z_\theta(s, a), \mathfrak{T}^\pi Z_\theta(s, a))$  for each  $s, a$  pairs. In terms of the representation  
 145 for  $Z_\theta(s, a)$ , we employ the unrestricted statistics, i.e., deterministic samples, due to its superiority in  
 146 MMDDL [17], instead of using predefined statistic functionals, e.g., quantiles in QR-DQN [8] or  
 147 histogram partitions in C51 [2]. More concretely, we use neural networks to generate samples that  
 148 approximate the value distribution. This can be expressed as  $Z_\theta(s, a) := \{Z_\theta(s, a)_i\}_{i=1}^N$ , where  $N$   
 149 is the number of generated samples. We refer to the samples  $\{Z_\theta(s, a)_i\}_{i=1}^N$  as *particles*. Then we  
 150 leverage the Dirac mixture  $\frac{1}{N} \sum_{i=1}^N \delta_{Z_\theta(s, a)_i}$  to approximate the true density function of  $Z^\pi(s, a)$ ,  
 151 thus minimizing the Sinkhorn divergence between the approximate distribution and its distributional  
 152 Bellman target. A detailed and generic distributional RL algorithm with Sinkhorn divergence and  
 153 particle representation is provided in Algorithm 1.

---

**Algorithm 1** Generic Sinkhorn distributional RL Update

---

**Require:** Number of generated samples  $N$ , the cost function  $c$  and hyperparameter  $\varepsilon$ .

**Input:** Sample transition  $(s, a, r', s')$

- 1: **if** Policy evaluation **then**
- 2:    $a^* \sim \pi(\cdot|s')$ .
- 3: **else**
- 4:    $a^* \leftarrow \arg \max_{a' \in \mathcal{A}} \frac{1}{N} \sum_{i=1}^N Z_\theta(s', a')_i$
- 5: **end if**
- 6:  $\mathfrak{T}Z_i \leftarrow r + \gamma Z_{\theta^*}(s', a^*)_i, \forall 1 \leq i \leq N$

**Output:**  $\overline{W}_{c,\varepsilon} \left( \{Z_\theta(s, a)_i\}_{i=1}^N, \{\mathfrak{T}Z_\theta(s, a)_j\}_{j=1}^N \right)$

---

154 **Remark.** By comparing the state-of-the-art MMDDRL algorithm [17], our Sinkhorn distributional  
155 RL simply modifies the distribution divergence. Hence, we can also easily extend our generic  
156 Sinkhorn algorithm to DQN-like architecture as well as IQN [7] and FQF [24]. A following question  
157 is whether there is any theoretical connection between Sinkhorn distributional RL and algorithms  
158 based on MMD and Wasserstein distance. We provide this crucial analysis in Section 4.2

## 159 4.2 Theoretical Analysis under Sinkhorn Divergence

160 **Convergence Analysis.** Firstly, we denote the supreme form of Sinkhorn divergence as  $\overline{W}_{c,\varepsilon}^\infty(\mu, \nu)$ :

161 
$$\overline{W}_{c,\varepsilon}^\infty(\mu, \nu) = \sup_{(x,a) \in \mathcal{S} \times \mathcal{A}} \overline{W}_{c,\varepsilon}(\mu(x, a), \nu(x, a)). \quad (10)$$

162 We will use  $\overline{W}_{c,\varepsilon}^\infty(\mu, \nu)$  to establish the convergence of  $\mathfrak{T}^\pi$  in Theorem 1.

163 **Theorem 1.** *If we leverage Sinkhorn loss  $\overline{W}_{c,\varepsilon}(\mu, \nu)$  in Eq. 8 as the distribution divergence in*  
164 *distributional RL, and **choose the unrectified kernel**  $k_\alpha := -\|x - y\|^\alpha$  as  $-c$  ( $\alpha > 0$ ), it holds that*

- 165 (1) *As  $\varepsilon \rightarrow 0$ ,  $\overline{W}_{c,\varepsilon}(\mu, \nu) \rightarrow 2W_\alpha(\mu, \nu)$ . When  $\varepsilon = 0$ ,  $\mathfrak{T}^\pi$  is a  $\gamma$ -contraction under  $\overline{W}_{c,\varepsilon}^\infty$ .*
- 166 (2) *As  $\varepsilon \rightarrow +\infty$ ,  $\overline{W}_{c,\varepsilon}(\mu, \nu) \rightarrow \text{MMD}_{k_\alpha}^2(\mu, \nu)$ . When  $\varepsilon = +\infty$ ,  $\mathfrak{T}^\pi$  is  $\gamma^{\alpha/2}$ -contractive under  $\overline{W}_{c,\varepsilon}^\infty$ .*
- 167 (3) *For any  $\varepsilon \in (0, +\infty)$ ,  $\mathfrak{T}^\pi$  is a **closely non-expansive operator** under  $\overline{W}_{c,\varepsilon}^\infty$ , and the difference*  
168 *term  $\Delta(\gamma) \rightarrow 0$  as  $\gamma \rightarrow 1$ .*

169 Proof is provided in Appendix B. Theorem 1 (1) and (2) are follow-up conclusions in terms of the  
170 convergence behavior of  $\mathfrak{T}^\pi$  based on the interpolation relationship between Sinkhorn divergence with  
171 Wasserstein distance and MMD [14]. Our key theoretical contribution is for the general  $\varepsilon \in (0, \infty)$ ,  
172 the convergence behavior is determined by the “joint” KL divergence in Eq. 9 between the optimal  
173 joint distribution  $\Pi^*$  and the Gibbs distribution associated with the cost function  $c$ . We conclude  
174 that  $\mathfrak{T}^\pi$  is a **close** non-expansive operator and the different term  $\Delta(\gamma) \rightarrow 0$  as  $\gamma \rightarrow 1$ . Note that  $\gamma$  is  
175 normally very close to 1 in practice, and this is beneficial for the convergence of  $\mathfrak{T}^\pi$  under  $\overline{W}_{c,\varepsilon}^\infty$ .

176 **Remark on Theorem 1 (3).** If we consider to use Gaussian kernel, we can not guarantee  $\mathfrak{T}^\pi$  is  
177 closely non-expansive for any  $\varepsilon \in (0, \infty)$ . This conclusion is consistent with those discussed  
178 in MMDDRL [17], where  $\mathfrak{T}^\pi$  is generally not a contraction operator under MMD equipped with  
179 Gaussian kernels as a counterexample has been pointed out in MMDDRL (when  $\varepsilon \rightarrow +\infty$ ). When  
180  $\varepsilon \rightarrow 0$ , the  $\gamma$ -contractive  $\mathfrak{T}^\pi$  under Wasserstein distance is also not contradictory to Theorem 1  
181 (3). Moreover, although we can only obtain that  $\mathfrak{T}^\pi$  is closely non-expansive, the expectation of  
182  $Z^\pi$  remains a  $\gamma$ -contraction (see Appendix B). In experiments, we thereby use  $k_\alpha$  and we can also  
183 demonstrate the appealing empirical performance of our SinkhornDRL algorithm in Section 5.

184 **Regularized Moment Matching under Sinkhorn Divergence.** We further examine the potential  
185 reason behind the empirical success for SinkhornDRL, although only a non-expansive contraction  
186 can be guaranteed for the general case when  $\varepsilon \in (0, +\infty)$  as shown in Theorem 1. Inspired by the  
187 similar manner in MMDDRL [17], we find that the Sinkhorn divergence with the Gaussian kernel  
188 can also promote to match all moments between two distributions. More specifically, the Sinkhorn  
189 divergence can be rewritten as a regularized moment matching form in Proposition 1.

190 **Proposition 1.** For  $\varepsilon \in (0, +\infty)$ , Sinkhorn divergence  $\overline{W}_{c,\varepsilon}(\mu, \nu)$  associated with Gaussian kernels  
 191  $k(x, y) = \exp(-(x - y)^2/(2\sigma^2))$  as  $-c$ , can be equivalent to

$$\overline{W}_{c,\varepsilon}(\mu, \nu) := \sum_{n=0}^{\infty} \frac{1}{\sigma^{2n} n!} \left( \tilde{M}_n(\mu) - \tilde{M}_n(\nu) \right)^2 + \varepsilon \mathbb{E} \left[ \log \frac{(\Pi_{\varepsilon}^*(X, Y))^2}{\Pi_{\varepsilon}^*(X, X') \Pi_{\varepsilon}^*(Y, Y')} \right], \quad (11)$$

192 where  $\Pi_{\varepsilon}^*$  denotes the optimal  $\Pi$  determined by  $\varepsilon$  by evaluating the Sinkhorn divergence via  
 193  $\min_{\Pi \in \Pi(\mu, \nu)} \overline{W}_{c,\varepsilon}(\mu, \nu)$ .  $\tilde{M}_n(\mu) = \mathbb{E}_{x \sim \mu} \left[ e^{-x^2/(2\sigma^2)} x^n \right]$ , and similarly for  $\tilde{M}_n(\nu)$ .

194 We provide the proof of Proposition 1 in Appendix C. Similar to MMDDRL associated with a  
 195 Gaussian kernel [17], Sinkhorn divergence approximately performs a regularized moment matching  
 196 scaled by  $e^{-x^2/(2\sigma^2)}$ . This similar moment matching impact intuitively explains the empirical success  
 197 of SinkhornDRL as MMDDRL, although the contraction of both MMD with Gaussian kernel [17]  
 198 and Sinkhorn divergence for general  $\varepsilon \in (0, +\infty)$  may not be guaranteed.

199 **Equivalence to Regularized MMD distributional RL.** Based on Proposition 1, we can immedi-  
 200 ately establish the connection between Sinkhorn divergence and MMD in Corollary 1, indicating that  
 201 minimizing Sinkhorn divergence between two distributions is equivalent to minimizing a regularized  
 202 squared MMD.

203 **Corollary 1.** For  $\varepsilon \in (0, +\infty)$  and denote  $\Pi_{\varepsilon}^*$  as the optimal  $\Pi$  by evaluating the Sinkhorn divergence,  
 204 it holds that

$$\overline{W}_{c,\varepsilon} := \text{MMD}_{-c}^2(\mu, \nu) + \varepsilon \mathbb{E} \left[ \log \frac{(\Pi_{\varepsilon}^*(X, Y))^2}{\Pi_{\varepsilon}^*(X, X') \Pi_{\varepsilon}^*(Y, Y')} \right], \quad (12)$$

205 where we use  $\overline{W}_{c,\varepsilon}$  to replace  $\overline{W}_{c,\varepsilon}(\mu, \nu)$  for short.

206 Proof of Corollary 1 is provided in Appendix C. It is worthy of noting that this equivalence is  
 207 established for the general case when  $\varepsilon \in (0, +\infty)$ , and it does not hold in the limit cases when  
 208  $\varepsilon \rightarrow 0$  or  $+\infty$ . For example, when  $\varepsilon \rightarrow +\infty$ , the second part including  $\varepsilon$  in Eq. 12 is not expected to  
 209 dominate. This is owing to the fact that the regularization term would be 0 as  $\Pi_{\varepsilon}^* \rightarrow \mu \otimes \nu$  when  
 210  $\varepsilon \rightarrow +\infty$ . In summary, even though the Sinkhorn divergence was initially proposed to serve as an  
 211 entropy regularized Wasserstein distance, it turns out that it is equivalent to a regularized MMD, as  
 212 revealed in Corollary 1. This connection provides strong evidence for our empirical results, in which  
 213 SinkhornDRL achieves competitive performance as opposed to MMDDRL.

### 214 4.3 Distributional RL via Sinkhorn Iterations

215 The theoretical analysis in Section 4.2 sheds light on the behavior of distributional RL with Sinkhorn  
 216 divergence, but another crucial issue we need to address is how to evaluate the Sinkhorn loss  
 217 effectively. Due to the advantages of Sinkhorn divergence that both enjoys geometry property of  
 218 optimal transport and the computational effectiveness of MMD, we can utilize Sinkhorn's algorithm,  
 219 i.e., Sinkhorn Iterations [21, 14], to evaluate the Sinkhorn loss. Notably, Sinkhorn iteration with  
 220  $L$  steps yields a differentiable and solvable efficiently loss function as the main burden involved  
 221 in it is the matrix-vector multiplication, which streams well on the GPU with simply adding extra  
 222 differentiable layers on the typical deep neural network, such as a DQN architecture.

223 Specifically, given two sample sequences  $\{Z_i\}_{i=1}^N, \{\mathfrak{Z}_j\}_{j=1}^N$  in the distributional RL algorithm, the  
 224 optimal transport distance is equivalent to the form:

$$\min_{P \in \mathbb{R}_+^{N \times N}} \left\{ \langle P, \hat{c} \rangle; P \mathbf{1}_N = \mathbf{1}_N, P^{\top} \mathbf{1}_N = \mathbf{1}_N \right\}, \quad (13)$$

225 where the empirical cost function  $\hat{c}_{i,j} = c(Z_i, \mathfrak{Z}_j)$ . By adding entropic regularization on optimal  
 226 transport distance, Sinkhorn divergence can be viewed to restrict the search space of  $P$  in the  
 227 following scaling form:

$$P_{i,j} = a_i \mathcal{K}_{i,j} b_j, \quad (14)$$

228 where  $\mathcal{K}_{i,j} = e^{-\hat{c}_{i,j}/\varepsilon}$  is the Gibbs kernel defined in Eq. 9. This allows us to leverage iterations  
 229 regarding the vectors  $a$  and  $b$ . More specifically, we initialize  $b_0 = \mathbf{1}_N$ , and then the Sinkhorn  
 230 iterations are expressed as

$$a_{l+1} \leftarrow \frac{\mathbf{1}_N}{\mathcal{K} b_l} \quad \text{and} \quad b_{l+1} \leftarrow \frac{\mathbf{1}_N}{\mathcal{K}^{\top} a_{l+1}}, \quad (15)$$

---

**Algorithm 2** Sinkhorn Iterations to Approximate  $\overline{\mathcal{W}}_{c,\varepsilon} \left( \{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N \right)$ 

---

**Input:** Two samples sequences  $\{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N$ , number of Sinkhorn iterations  $L$  and hyperparameter  $\varepsilon$ .

- 1:  $\hat{c}_{i,j} = c(Z_i, \mathfrak{T}Z_j)$  for  $\forall i = 1, \dots, N, j = 1, \dots, N$
- 2:  $\mathcal{K}_{i,j} = \exp(-\hat{c}_{i,j}/\varepsilon)$
- 3:  $b_0 \leftarrow \mathbf{1}_N$
- 4: **for**  $l = 1, 2, \dots, L$  **do**
- 5:      $a_l \leftarrow \frac{\mathbf{1}_N}{\mathcal{K}b_{l-1}}, b_l \leftarrow \frac{\mathbf{1}_N}{\mathcal{K}a_l}$
- 6: **end for**
- 7:  $\widehat{\mathcal{W}}_{c,\varepsilon} \left( \{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N \right) = \langle (K \odot \hat{c})b, a \rangle$

**Return:**  $\widehat{\mathcal{W}}_{c,\varepsilon} \left( \{Z_i\}_{i=1}^N, \{\mathfrak{T}Z_j\}_{j=1}^N \right)$ 

---

231 where  $\dot{\cdot}$  indicates an entry-wise division. It has been proven that Sinkhorn iteration asymptotically  
232 converges to the true loss in a linear rate [14, 12, 6]. We provide a detailed algorithm description of  
233 Sinkhorn iterations in Algorithm 2. With the efficient and differential Sinkhorn iterations, we can  
234 easily evaluate the Sinkhorn divergence and thus let our algorithm enjoy its theoretical advantages. In  
235 practice, we need to choose  $L$  and  $\varepsilon$ , and we conduct a rigorous sensitivity analysis in Section 5.

## 236 5 Experiments

237 We demonstrate the effectiveness of SinkhornDRL as described in Algorithm 1 on the full 55 Atari  
238 2600 games. Specifically, we leverage the same architecture as QR-DQN [8], and replace the quantiles  
239 output with  $N$  particles, i.e., samples. In contrast to MMDDRL, SinkhornDRL only changes the  
240 distribution divergence from MMD to Sinkhorn divergence, and therefore the potential superiority in  
241 the performance can be attributed to the advantages of Sinkhorn divergence. In Section 5.1, we make  
242 a rigorous comparison between SinkhornDRL with other typical distributional RL algorithms from  
243 the perspectives of learning curves and final ratio improvement of returns. An extensive sensitivity  
244 analysis in terms of multiple hyperparameters in SinkhornDRL is provided in Section 5.2.

245 **Baselines.** Due to the interpolation characteristic of Sinkhorn divergence between Wasserstein  
246 distance and MMDDRL, we choose three typical distributional RL algorithms as classic baselines,  
247 including QR-DQN [8] that approximates the Wasserstein distance, C51 [2] and MMDDRL [17], as  
248 well as DQN [16]. MMDDRL algorithm is implemented with the same architecture as QRDQN, and  
249 leverages Gaussian kernels  $k_h(x, y) = \exp(-(x - y)^2/h)$  with the kernel mixture trick covering a  
250 range of bandwidths  $h$ , which is same as the basic setting in the original MMDDQN paper [17]. We  
251 deploy all algorithms on 55 Atari 2600 games, and reported results are averaged over 3 seeds with  
252 the shade indicating the standard deviation.

253 **Hyperparameter settings.** For a fair comparison with QR-DQN, C51 and MMDDRL, we used  
254 the same hyperparameters: the number of generated samples  $N = 200$ , Adam optimizer with  
255  $\text{lr} = 0.00005, \epsilon_{\text{Adam}} = 0.01/32$ . We used a target network to compute the distributional Bellman  
256 target, which fits well in the neural Z-fitted iteration framework. In addition, we choose number of  
257 Sinkhorn iterations  $L = 10$  and smoothing hyperparameter  $\varepsilon = 10.0$  in Section 5.1 as they are not  
258 sensitive within a proper interval as demonstrated in Section 5.2. We choose the unrectified kernel as  
259 the cost function, i.e.,  $-c = k_\alpha$ , and select  $\alpha = 2$  in  $k_\alpha$  in our SinkhornDRL algorithm.

### 260 5.1 Performance of SinkhornDRL

261 Figure 1 illustrates that SinkhornDRL can achieve the competitive performance across 55 Atari games  
262 compared with various baseline algorithms with different metrics  $d_p$  and representation manners on  
263  $Z_\theta$ . On a large number of games, e.g., Tennis, Seaquest and Atlantis, SinkhornDRL can significantly  
264 outperform other baselines, especially on Tennis where other algorithms even fail to converge. The  
265 improvement of SinkhornDRL over MMDDRL empirically verifies the regularization advantage of  
266 the Sinkhorn as analyzed in Corollary 1. On some games, e.g., Breakout, Pong and SpaceInvaders,

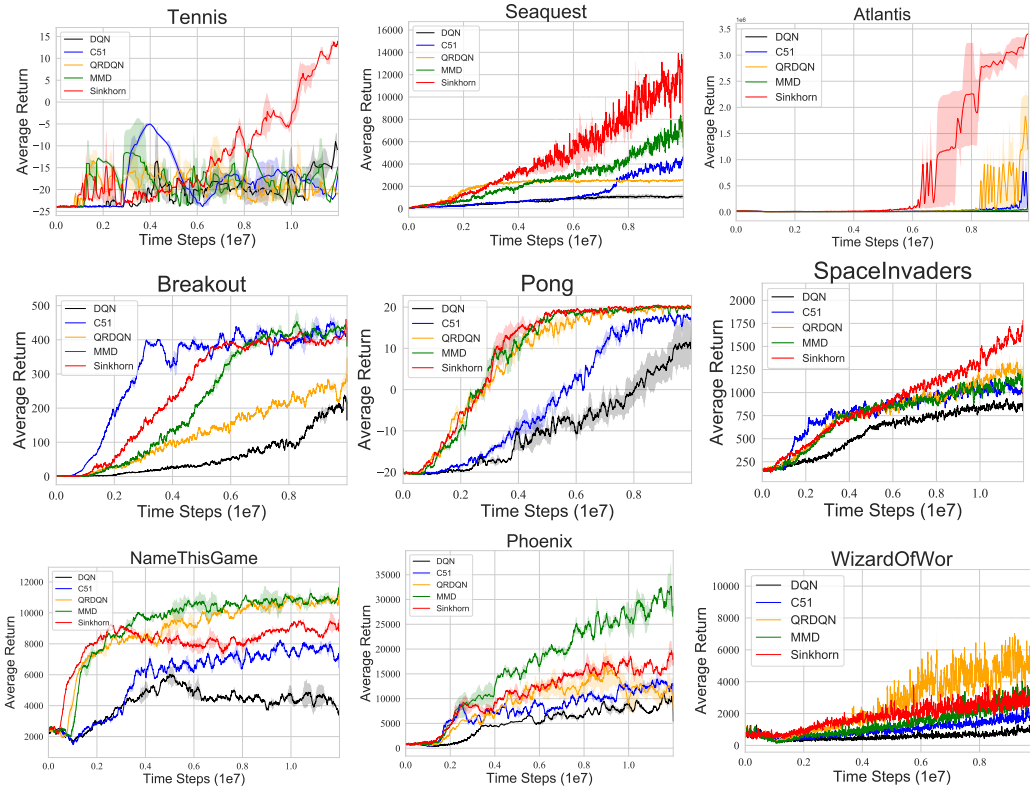


Figure 1: Learning curves of SinkhornDRL algorithm compared with DQN, C51, QR-DQN and MMD, on nine typical Atari games over 3 seeds.

267 SinkhornDRL is on par with MMDDRL and other baselines, while on the last row in Figure 1,  
 268 SinkhornDRL is slightly inferior to the state-of-the-art algorithm. We provide learning curves of all  
 269 typical distributional RL algorithms on all 55 Atari games in Appendix E, where SinkhornDRL still  
 270 achieves the competitive performance in general.

271 To further demonstrate theoretical properties of SinkhornDRL in Theorem 1, we conduct a ratio im-  
 272 provement comparison across 55 Atari games between SinkhornDRL with QR-DQN and MMDDRL,  
 273 respectively. Figure 2 showcases that by comparing with QR-DQN (left), SinkhornDRL achieves  
 274 better performance across more than half of considered games. More importantly, the superiority  
 275 of SinkhornDRL is significant across a large amount of games, including Venture, Seaquest, Tennis  
 276 and Phoenix. This empirical outperformance verifies the effectiveness and potential of smoothing  
 277 Wasserstein distance in distributional RL, e.g., Sinkhorn divergence. In contrast with MMDDRL, the  
 278 superiority of SinkhornDRL is reduced with the performance improvement only on a small proportion  
 279 of games, while a remarkable boost of performance for SinkhornDRL on a large amount of games  
 280 can be easily observed. We also report mean and median of best human-normalized scores in Table 2  
 281 of Appendix D, where SinkhornDRL achieves almost state-of-the-art performance as MMDDRL on  
 282 average.

283 Therefore, we conclude that SinkhornDRL is competitive with the state-of-the-art distributional  
 284 RL algorithms, e.g., MMDDRL, and can be extremely superior over existing algorithms on a large  
 285 proportion of games. This empirical success can be owing to theoretical advantage of Sinkhorn  
 286 divergence that simultaneously makes full use of the data geometry from Wasserstein distance and  
 287 the unbiased gradient estimate property from MMD, which coincides with results in Theorem 1.

## 288 5.2 Sensitivity Analysis and Computational Cost

289 Figure 3 (a) suggests the performance of our algorithm is robust to  $\varepsilon$  in a certain range, e.g., [1, 500],  
 290 facilitating its deployment in practice. If we increase  $\varepsilon$ , SinkhornDRL’s performance tends to MMD,

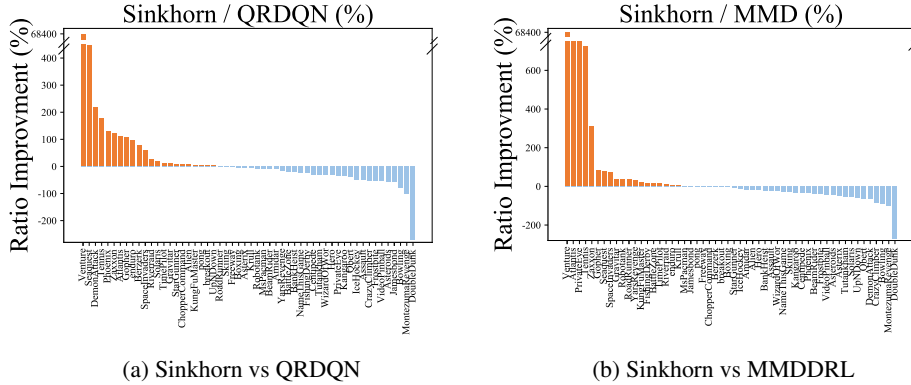


Figure 2: Ratio improvement of return for Sinkhorn distributional RL algorithm over QRDQN (left) and MMDDRL (right) over 3 seeds. For example, the ratio improvement is calculated by  $(\text{Sinkhorn} - \text{QRDQN}) / \text{QRDQN}$  in the left.

291 while if we gradually decline  $\varepsilon$ , SinkhornDRL’s performance tends to QR-DQN. It is also noted that  
 292 Sinkhorn iterations in Algorithm 2 will suffer from the numerical instability issue under an overly  
 293 small or large  $\varepsilon$ . More results with the discussion are provided in Appendix F. It is also illustrated  
 294 that our algorithm is insensitive to the number of iterations  $L$  and samples  $N$  as well, but an overly  
 295 large  $N$  can slightly worsen the performance of SinkhornDRL, and at the same time increases the  
 296 computational burden. Therefore, a proper number of samples, e.g., 200, is sufficient to attain an  
 297 appealing performance with the computational effectiveness.

298 For the computation cost, SinkhornDRL indeed increases around 50% computation cost compared  
 299 with QR-DQN and C51, but only slightly increases the cost (by around 20%) in contrast to MMDDRL.  
 300 Detailed comparison is given in Appendix F.

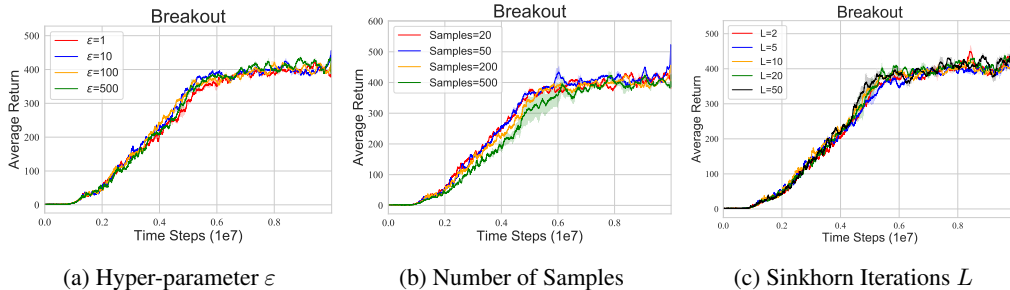


Figure 3: Sensitivity analysis of SinkhornDRL on Breakout regarding  $\varepsilon$ , number of samples, and number of iteration  $L$ . Learning curves are reported over 3 seeds.

## 301 6 Discussions and Conclusion

302 The main limitation of our proposal is that the superiority over existing state-of-the-art algorithms may  
 303 not be sufficiently significant. To extend our algorithm for better performance, implicit generative  
 304 models, including parameterizing the cost function in Sinkhorn loss, can be further incorporated. We  
 305 leave it as the future work. Moreover, other divergences, e.g., those that can also smooth Wasserstein  
 306 distance, can also be applied into the design of distributional RL algorithms in the future.

307 In this paper, a novel family of distributional RL algorithms based on Sinkhorn Divergence is proposed  
 308 that accomplishes a competitive performance compared with the-state-of-the-art distributional RL  
 309 algorithms on 55 Atari games. Theoretical analysis about the convergence and moment matching  
 310 behavior is provided along with a rigorous empirical verification. Albeit being associated with MMD  
 311 algorithm, distributional RL with Sinkhorn divergence is complementary to previous algorithms,  
 312 leading to an important contribution among the research community.

## 313 References

- 314 [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial  
315 networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- 316 [2] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforce-  
317 ment learning. *International Conference on Machine Learning (ICML)*, 2017.
- 318 [3] Marc G Bellemare, Ivo Danihelka, Will Dabney, Shakir Mohamed, Balaji Lakshminarayanan,  
319 Stephan Hoyer, and Rémi Munos. The cramer distance as a solution to biased wasserstein  
320 gradients. *arXiv preprint arXiv:1705.10743*, 2017.
- 321 [4] Pablo Samuel Castro, Subhodeep Moitra, Carles Gelada, Saurabh Kumar, and Marc G Bellemare.  
322 Dopamine: A research framework for deep reinforcement learning. *CoRR abs/1812.06110*,  
323 2018.
- 324 [5] Yutian Chen, Max Welling, and Alex Smola. Super-samples from kernel herding. *UAI, 109–116*.  
325 *AUAI Press*, 2012.
- 326 [6] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in*  
327 *neural information processing systems*, 26, 2013.
- 328 [7] Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for  
329 distributional reinforcement learning. *International Conference on Machine Learning (ICML)*,  
330 2018.
- 331 [8] Will Dabney, Mark Rowland, Marc G Bellemare, and Rémi Munos. Distributional reinforcement  
332 learning with quantile regression. *Association for the Advancement of Artificial Intelligence*  
333 *(AAAI)*, 2018.
- 334 [9] Odin Elie and Charpentier Arthur. *Dynamic Programming in Distributional Reinforcement*  
335 *Learning*. PhD thesis, Université du Québec à Montréal, 2020.
- 336 [10] Jianqing Fan, Zhaoran Wang, Yuchen Xie, and Zhuoran Yang. A theoretical analysis of deep  
337 q-learning. In *Learning for Dynamics and Control*, pages 486–489. PMLR, 2020.
- 338 [11] Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trouvé, and  
339 Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences.  
340 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2681–2690.  
341 PMLR, 2019.
- 342 [12] Joel Franklin and Jens Lorenz. On the scaling of multidimensional matrices. *Linear Algebra*  
343 *and its applications*, 114:717–735, 1989.
- 344 [13] Aude Genevay, Lénaïc Chizat, Francis Bach, Marco Cuturi, and Gabriel Peyré. Sample  
345 complexity of sinkhorn divergences. In *The 22nd International Conference on Artificial*  
346 *Intelligence and Statistics*, pages 1574–1583. PMLR, 2019.
- 347 [14] Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn  
348 divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–  
349 1617. PMLR, 2018.
- 350 [15] Christian Léonard. A survey of the schrödinger problem and some of its connections with  
351 optimal transport. *arXiv preprint arXiv:1308.0215*, 2013.
- 352 [16] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
353 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
354 Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 355 [17] Thanh Tang Nguyen, Sunil Gupta, and Svetha Venkatesh. Distributional reinforcement learning  
356 with maximum mean discrepancy. *Association for the Advancement of Artificial Intelligence*  
357 *(AAAI)*, 2020.
- 358 [18] Aaditya Ramdas, Nicolás García Trillos, and Marco Cuturi. On wasserstein two-sample testing  
359 and related families of nonparametric tests. *Entropy*, 19(2):47, 2017.



- 360 [19] Mark Rowland, Marc Bellemare, Will Dabney, Rémi Munos, and Yee Whye Teh. An analysis  
361 of categorical distributional reinforcement learning. In *International Conference on Artificial*  
362 *Intelligence and Statistics*, pages 29–37. PMLR, 2018.
- 363 [20] Ludger Rüschendorf and Wolfgang Thomsen. Closedness of sum spaces and the generalized  
364 schrödinger problem. *Theory of Probability & Its Applications*, 42(3):483–494, 1998.
- 365 [21] Richard Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums.  
366 *The American Mathematical Monthly*, 74(4):402–405, 1967.
- 367 [22] Gábor J Székely. E-statistics: The energy of statistical samples. *Bowling Green State University,*  
368 *Department of Mathematics and Statistics Technical Report*, 3(05):1–18, 2003.
- 369 [23] Eric Wong, Frank Schmidt, and Zico Kolter. Wasserstein adversarial examples via projected  
370 sinkhorn iterations. In *International Conference on Machine Learning*, pages 6808–6817.  
371 PMLR, 2019.
- 372 [24] Derek Yang, Li Zhao, Zichuan Lin, Tao Qin, Jiang Bian, and Tie-Yan Liu. Fully parameterized  
373 quantile function for distributional reinforcement learning. *Advances in neural information*  
374 *processing systems*, 32:6193–6202, 2019.
- 375 [25] Shangtong Zhang. Modularized implementation of deep rl algorithms in pytorch. <https://github.com/ShangtongZhang/DeepRL>, 2018.  
376
- 377 [26] Fan Zhou, Jianing Wang, and Xingdong Feng. Non-crossing quantile regression for distri-  
378 butional reinforcement learning. *Advances in Neural Information Processing Systems*, 33,  
379 2020.
- 380 [27] Florian Ziel. The energy distance for ensemble and scenario reduction. *arXiv preprint*  
381 *arXiv:2005.14670*, 2020.

## 382 Checklist

- 383 1. For all authors...
- 384 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
385 contributions and scope? [Yes]
- 386 (b) Did you describe the limitations of your work? [Yes] We provide the discussion about  
387 the limitation of our proposal in Section 6.
- 388 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 389 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
390 them? [Yes]
- 391 2. If you are including theoretical results...
- 392 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Please refer to  
393 Appendix B and C.
- 394 (b) Did you include complete proofs of all theoretical results? [Yes] Please refer to  
395 Appendix B and C.
- 396 3. If you ran experiments...
- 397 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
398 mental results (either in the supplemental material or as a URL)? [Yes]
- 399 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
400 were chosen)? [Yes] Our implementation is adapted from Pytorch distributional RL  
401 modules [25].
- 402 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
403 ments multiple times)? [Yes]
- 404 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
405 of GPUs, internal cluster, or cloud provider)? [Yes] We provide the comparison of  
406 computational cost in Figure 12 of Appendix F.

- 407 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 408 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 409 (b) Did you mention the license of the assets? [N/A]
- 410 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 411
- 412 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 413 using/curating? [N/A]
- 414 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 415 information or offensive content? [N/A]
- 416 5. If you used crowdsourcing or conducted research with human subjects...
- 417 (a) Did you include the full text of instructions given to participants and screenshots, if
- 418 applicable? [N/A]
- 419 (b) Did you describe any potential participant risks, with links to Institutional Review
- 420 Board (IRB) approvals, if applicable? [N/A]
- 421 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 422 spent on participant compensation? [N/A]

## 423 A Definition of distances and Contraction

424 **Definition of distances.** Given two random variables  $X$  and  $Y$ ,  $p$ -Wasserstein metric  $W_p$  between  
425 the distributions of  $X$  and  $Y$  is defined as

$$W_p(X, Y) = \left( \int_0^1 |F_X^{-1}(\omega) - F_Y^{-1}(\omega)|^p d\omega \right)^{1/p} = \|F_X^{-1} - F_Y^{-1}\|_p, \quad (16)$$

426 which  $F^{-1}$  is the inverse cumulative distribution function of a random variable with the cumulative  
427 distribution function as  $F$ . Further,  $\ell_p$  distance [9] is defined as

$$\ell_p(X, Y) := \left( \int_{-\infty}^{\infty} |F_X(\omega) - F_Y(\omega)|^p d\omega \right)^{1/p} = \|F_X - F_Y\|_p \quad (17)$$

428 The  $\ell_p$  distance and Wasserstein metric are identical at  $p = 1$ , but are otherwise distinct. Note that  
429 when  $p = 2$ ,  $\ell_p$  distance is also called Cramér distance [3]  $d_C(X, Y)$ . Also, the Cramér distance has  
430 a different representation given by

$$d_C(X, Y) = \mathbb{E}|X - Y| - \frac{1}{2}\mathbb{E}|X - X'| - \frac{1}{2}\mathbb{E}|Y - Y'|, \quad (18)$$

431 where  $X'$  and  $Y'$  are the i.i.d. copies of  $X$  and  $Y$ . Energy distance [22, 27] is a natural extension of  
432 Cramér distance to the multivariate case, which is defined as

$$d_E(\mathbf{X}, \mathbf{Y}) = \mathbb{E}\|\mathbf{X} - \mathbf{Y}\| - \frac{1}{2}\mathbb{E}\|\mathbf{X} - \mathbf{X}'\| - \frac{1}{2}\mathbb{E}\|\mathbf{Y} - \mathbf{Y}'\|, \quad (19)$$

433 where  $\mathbf{X}$  and  $\mathbf{Y}$  are multivariate. Moreover, the energy distance is a special case of the maximum  
434 mean discrepancy (MMD), which is formulated as

$$\text{MMD}(\mathbf{X}, \mathbf{Y}; k) = (\mathbb{E}[k(\mathbf{X}, \mathbf{X}')] + \mathbb{E}[k(\mathbf{Y}, \mathbf{Y}')] - 2\mathbb{E}[k(\mathbf{X}, \mathbf{Y})])^{1/2} \quad (20)$$

435 where  $k(\cdot, \cdot)$  is a continuous kernel on  $\mathcal{X}$ . In particular, if  $k$  is a trivial kernel, MMD degenerates  
436 to energy distance. Additionally, we further define the supreme MMD, which is a functional  
437  $\mathcal{P}(\mathcal{X})^{\mathcal{S} \times \mathcal{A}} \times \mathcal{P}(\mathcal{X})^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}$  defined as

$$\text{MMD}_\infty(\mu, \nu) = \sup_{(x, a) \in \mathcal{S} \times \mathcal{A}} \text{MMD}_\infty(\mu(x, a), \nu(x, a)) \quad (21)$$

438 We further present the convergence rate under different distribution divergences.

- 439 •  $\mathcal{T}^\pi$  is  $\gamma$ -contractive under the supreme form of Wasserstein distance  $W_p$ .
- 440 •  $\mathcal{T}^\pi$  is  $\gamma^{1/p}$ -contractive under the supreme form of  $\ell_p$  distance.
- 441 •  $\mathcal{T}^\pi$  is  $\gamma^{\alpha/2}$ -contractive under  $\text{MMD}_\infty$  with the kernel  $k_\alpha(x, y) = -\|x - y\|^\alpha, \forall \alpha > 0$ .

### 442 Proof of Contraction.

- 443 • Contraction under supreme form of Wasserstein distance is provided in Lemma 3 [2].
- 444 • Contraction under supreme form of  $\ell_p$  distance can refer to Theorem 3.4 [9].
- 445 • Contraction under  $\text{MMD}_\infty$  is provided in Lemma 6 [17].

## 446 B Proof of Theorem 1

447 *Proof. 1.* As  $\varepsilon \rightarrow 0$  and  $c = -k_\alpha$ , it is obvious to observe that Sinkhorn loss degenerates to the  
448 Wasserstein distance. We also have the conclusion that the distributional Bellman operator  $\mathfrak{T}^\pi$  is  
449  $\gamma$ -contractive under the supreme form of Wasserstein distance, the proof of which is provided in  
450 Lemma 3 [2]. Since the above conclusion is made directly based on the limiting case when  $\varepsilon = 0$ , for  
451 an unspecified  $\varepsilon$ , we need a more rigorous proof. We show that their distance difference is **at most**  
452 **an infinitesimal**  $\delta$ .

453 Firstly, as  $\mathcal{W}_{c,\varepsilon} \rightarrow W_\alpha$  and the regularization term is non-negative, using the language of  $(\varepsilon, \delta)$   
 454 definition, we have: for  $\forall \delta$ , there exists a small positive constant  $a$ , such that  $\mathcal{W}_{c,\varepsilon} - W_\alpha < \delta$  when  
 455  $\varepsilon \leq a$ . Based on that, we have the contraction conclusion:

$$\begin{aligned} \overline{\mathcal{W}}_{-\kappa_\alpha,\varepsilon}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) &= \overline{\mathcal{W}}_{-\kappa_\alpha,\varepsilon}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) - W_\alpha^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) + W_\alpha^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) \\ &\leq \delta + W_\alpha^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2), \end{aligned} \quad (22)$$

456 where the second term  $W_\alpha^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2)$  is contractive, and thus for the unspecified  $\varepsilon$ , the only  
 457 difference from the limiting  $\varepsilon = 0$  is an infinitesimal  $\delta$ , which will vanish as  $\varepsilon \rightarrow 0$  or  $(a \rightarrow 0)$ .

2. As  $\varepsilon \rightarrow \infty$ , our complete proof is inspired by [18, 14]. Recap the definition of squared MMD is

$$\mathbb{E}[k(\mathbf{X}, \mathbf{X}')] + \mathbb{E}[k(\mathbf{Y}, \mathbf{Y}')] - 2\mathbb{E}[k(\mathbf{X}, \mathbf{Y})]$$

When the kernel function  $k$  degenerates to a unrectified  $k_\alpha(x, y) := -\|x - y\|^\alpha$  for  $\alpha \in (0, 2)$ , the squared MMD would degenerate to

$$\mathbb{E}\|\mathbf{X} - \mathbf{X}'\|^\alpha + \mathbb{E}\|\mathbf{Y} - \mathbf{Y}'\|^\alpha - 2\mathbb{E}\|\mathbf{X} - \mathbf{Y}\|^\alpha$$

On the other hand, we have the Sinkhorn loss as

$$\overline{\mathcal{W}}_{c,\infty}(\mu, \nu) = 2\mathcal{W}_{c,\infty}(\mu, \nu) - \mathcal{W}_{c,\infty}(\nu, \nu) - \mathcal{W}_{c,\infty}(\mu, \mu)$$

Denoting  $\Pi_\varepsilon$  be the unique minimizer for  $\overline{\mathcal{W}}_{c,\varepsilon}$ , it holds that  $\Pi_\varepsilon \rightarrow \mu \otimes \nu$  as  $\varepsilon \rightarrow \infty$ . That being  
 said,  $\mathcal{W}_{c,\infty}(\mu, \nu) \rightarrow \int c(x, y)d\mu(x)d\nu(y) + 0 = \int c(x, y)d\mu(x)d\nu(y)$ . If  $c = -k_\alpha = -\|x - y\|^\alpha$ ,  
 we eventually have  $\mathcal{W}_{-k_\alpha,\infty}(\mu, \nu) \rightarrow \int \|x - y\|^\alpha d\mu(x)d\nu(y) = \mathbb{E}\|\mathbf{X} - \mathbf{Y}\|^\alpha$ . Finally, we can have

$$\overline{\mathcal{W}}_{-k_\alpha,\infty} \rightarrow 2\mathbb{E}\|\mathbf{X} - \mathbf{Y}\|^\alpha - \mathbb{E}\|\mathbf{X} - \mathbf{X}'\|^\alpha - \mathbb{E}\|\mathbf{Y} - \mathbf{Y}'\|^\alpha$$

458 which is exactly the form of squared MMD. Now the key is prove that  $\Pi_\varepsilon \rightarrow \mu \otimes \nu$  as  $\varepsilon \rightarrow \infty$ .

Firstly, it is apparent that  $\mathcal{W}_{c,\varepsilon}(\mu, \nu) \leq \int c(x, y)d\mu(x)d\nu(y)$  as  $\mu \otimes \nu \in \Pi(\mu, \nu)$ . Let  $\{\varepsilon_k\}$  be a  
 positive sequence that diverges to  $\infty$ , and  $\Pi_k$  be the corresponding sequence of unique minimizers for  
 $\mathcal{W}_{c,\varepsilon}$ . According to the optimality condition, it must be the case that  $\int c(x, y)d\Pi_k + \varepsilon_k \text{KL}(\Pi_k, \mu \otimes \nu) \leq \int c(x, y)d\mu \otimes \nu + 0$  (when  $\Pi(\mu, \nu) = \mu \otimes \nu$ ). Thus,

$$\text{KL}(\Pi_k, \mu \otimes \nu) \leq \frac{1}{\varepsilon_k} \left( \int c d\mu \otimes \nu - \int c d\Pi_k \right) \rightarrow 0.$$

Besides, by the compactness of  $\Pi(\mu, \nu)$ , we can extract a converging subsequence  $\Pi_{n_k} \rightarrow \Pi_\infty$ .  
 Since KL is weakly lower-semicontinuous, it holds that

$$\text{KL}(\Pi_\infty, \mu \otimes \nu) \leq \liminf_{k \rightarrow \infty} \text{KL}(\Pi_{n_k}, \mu \otimes \nu) = 0$$

459 Hence  $\Pi_\infty = \mu \otimes \nu$ . That being said that the optimal coupling is simply the product of the marginals,  
 460 indicating that  $\Pi_\varepsilon \rightarrow \mu \otimes \nu$  as  $\varepsilon \rightarrow \infty$ . As a special case, when  $\alpha = 1$ ,  $\overline{\mathcal{W}}_{-k_1,\infty}(u, v)$  is equivalent  
 461 to the energy distance

$$d_E(\mathbf{X}, \mathbf{Y}) := 2\mathbb{E}\|\mathbf{X} - \mathbf{Y}\| - \mathbb{E}\|\mathbf{X} - \mathbf{X}'\| - \mathbb{E}\|\mathbf{Y} - \mathbf{Y}'\|. \quad (23)$$

462 In summary, if the cost function is the rectified kernel  $k_\alpha$ , it is the case that  $\overline{\mathcal{W}}_{-k_\alpha,\varepsilon}$  converges to the  
 463 squared MMD as  $\varepsilon \rightarrow \infty$ . According to [17],  $\mathfrak{T}^\pi$  is  $\gamma^{\alpha/2}$ -contractive in the supreme form of MMD  
 464 with the rectified kernel  $k_\alpha$ .

465 For the unspecified  $\varepsilon$ , we can get the similar result to the case of  $\varepsilon \rightarrow 0$ . For  $\forall \delta$ , there exists a large  
 466 positive constant  $M$ , such that  $\text{MMD}_{k_\alpha}^2 - \mathcal{W}_{c,\varepsilon} < \delta$  when  $\varepsilon \geq M$ . Based on that, we have the  
 467 contraction conclusion:

$$\begin{aligned} \overline{\mathcal{W}}_{-\kappa_\alpha,\varepsilon}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) &= \overline{\mathcal{W}}_{-\kappa_\alpha,\varepsilon}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) - \text{MMD}_\infty^2(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) + \text{MMD}_\infty^2(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) \\ &\leq \text{MMD}_\infty^2(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) - \delta, \end{aligned} \quad (24)$$

468 where the first term  $\text{MMD}_\infty^2(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2)$  is  $\gamma^{\frac{\alpha}{2}}$ -contractive, and thus for the unspecified  $\varepsilon$ , the  
 469 only difference from the limiting  $\varepsilon = \infty$  is an infinitesimal  $\delta$ , which will vanish as  $\varepsilon \rightarrow +\infty$  or  
 470  $(M \rightarrow +\infty)$ .

471 **3.** For  $\varepsilon \in (0, +\infty)$ , a key observation for the analysis is that the Sinkhorn divergence would  
 472 degenerate to a two-dimensional KL divergence, and therefore embraces a similar convergence  
 473 behavior to KL divergence. Concretely, according to the equivalent form of  $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$  in Eq. 9, it  
 474 can be expressed as the KL divergence between an optimal joint distribution and a Gibbs distribution  
 475 associated with the cost function:

$$\mathcal{W}_{c,\varepsilon}(\mu, \nu) := \text{KL}(\Pi^*(\mu, \nu) | \mathcal{K}(\mu, \nu)), \quad (25)$$

476 where  $\Pi^*$  is the optimal joint distribution. Thus, the total Sinkhorn divergence is expressed as

$$\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu) := 2\text{KL}(\Pi^*(\mu, \nu) | \mathcal{K}(\mu, \nu)) - \text{KL}(\Pi^*(\mu, \mu) | \mathcal{K}(\mu, \mu)) - \text{KL}(\Pi^*(\nu, \nu) | \mathcal{K}(\nu, \nu)). \quad (26)$$

477 Due to the form of  $\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu)$ , the convergence behavior is determined by  $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$ , which is  
 478 similar to the behavior of KL divergence. Thus, we will focus on the convergence analysis of  
 479  $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$ . We firstly elaborate a Lemma regarding to the convergence under KL divergence.

480 **Lemma 1.** Denote the supreme of  $D_{\text{KL}}$  as  $D_{\text{KL}}^\infty$ , we have: (1)  $\mathfrak{T}^\pi$  is a non-expansive operator under  
 481  $D_{\text{KL}}^\infty$ , i.e.,  $D_{\text{KL}}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) \leq D_{\text{KL}}^\infty(Z_1, Z_2)$ , (2) the expectation of  $Z^\pi$  is still  $\gamma$ -contractive under  
 482  $D_{\text{KL}}^\infty$ , i.e.,  $\|\mathbb{E}\mathfrak{T}^\pi Z_1 - \mathbb{E}\mathfrak{T}^\pi Z_2\|_\infty \leq \gamma \|\mathbb{E}Z_1 - \mathbb{E}Z_2\|_\infty$ .

483 *Proof.* (1) We recap three crucial properties of a divergence metric. The first is *scale sensitive (S)*  
 484 (of order  $\beta$ ,  $\beta > 0$ ), i.e.,  $d_p(cX, cY) \leq |c|^\beta d_p(X, Y)$ . The second property is *shift invariant (I)*,  
 485 i.e.,  $d_p(A + X, A + Y) \leq d_p(X, Y)$ . The last one is *unbiased gradient (U)*. We use  $p$  and  $q$  to  
 486 denote the density function of two random variables  $X$  and  $Y$ , and thus  $D_{\text{KL}}(X, Y)$  is defined as  
 487  $D_{\text{KL}}(X, Y) = \int_{-\infty}^{\infty} p(x) \frac{p(x)}{q(x)} dx$ . Firstly, we show that  $D_{\text{KL}}(X, Y)$  is NOT scale sensitive:

$$\begin{aligned} D_{\text{KL}}(aX, aY) &= \int_{-\infty}^{\infty} \frac{1}{a} p\left(\frac{x}{a}\right) \log \frac{\frac{1}{a} p\left(\frac{x}{a}\right)}{\frac{1}{a} q\left(\frac{x}{a}\right)} dx \\ &= \int_{-\infty}^{\infty} p(y) \log \frac{p(y)}{q(y)} dy \\ &= D_{\text{KL}}(X, Y), \text{ with } \beta = 0 \end{aligned} \quad (27)$$

488 We further show that  $D_{\text{KL}}(X, Y)$  is shift invariant:

$$\begin{aligned} D_{\text{KL}}(A + X, A + Y) &= \int_{-\infty}^{\infty} p(x - A) \log \frac{p(x - A)}{q(x - A)} dx \\ &= \int_{-\infty}^{\infty} p(y) \log \frac{p(y)}{q(y)} dy \\ &= D_{\text{KL}}(X, Y) \end{aligned} \quad (28)$$

489 Moreover, it is well-known that KL divergence has unbiased sample gradients [3]. The supreme  $D_{\text{KL}}$   
 490 is a functional  $\mathcal{P}(\mathcal{X})^{\mathcal{S} \times \mathcal{A}} \times \mathcal{P}(\mathcal{X})^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}$  defined as

$$D_{\text{KL}}^\infty(\mu, \nu) = \sup_{(x,a) \in \mathcal{S} \times \mathcal{A}} D_{\text{KL}}(\mu(x, a), \nu(x, a)) \quad (29)$$

491 Therefore, we prove  $\mathfrak{T}^\pi$  is at best a non-expansive operator under the supreme form of  $D_{\text{KL}}$ :

$$\begin{aligned} &D_{\text{KL}}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) \\ &= \sup_{s,a} D_{\text{KL}}(\mathfrak{T}^\pi Z_1(s, a), \mathfrak{T}^\pi Z_2(s, a)) \\ &= \sup_{s,a} D_{\text{KL}}(R(s, a) + \gamma Z_1(S', A'), R(s, a) + \gamma Z_2(S', A')) \\ &= D_{\text{KL}}(Z_1(S', A'), Z_2(S', A')) \\ &\leq \sup_{s',a'} D_{\text{KL}}(Z_1(s', a'), Z_2(s', a')) \\ &= D_{\text{KL}}^\infty(Z_1, Z_2) \end{aligned} \quad (30)$$

492 There we have  $D_{\text{KL}}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) \leq D_{\text{KL}}^\infty(Z_1, Z_2)$ , implying that  $\mathfrak{T}^\pi$  is a non-expansive operator  
 493 under  $D_{\text{KL}}^\infty$ .

494 (2) This statement is an immediate conclusion based on the Lemma 4 in [2]. We give the proof for  
 495 the completeness. This conclusion holds because the  $\mathfrak{T}^\pi$  degenerates to  $\mathcal{T}^\pi$  regardless of the metric  
 496  $d_p$ . Specifically, due to the linearity of expectation, we obtain that

$$\|\mathbb{E}\mathfrak{T}^\pi Z_1 - \mathbb{E}\mathfrak{T}^\pi Z_2\|_\infty = \|\mathcal{T}^\pi \mathbb{E}Z_1 - \mathcal{T}^\pi \mathbb{E}Z_2\|_\infty \leq \gamma \|\mathbb{E}Z_1 - \mathbb{E}Z_2\|_\infty. \quad (31)$$

497 This implies that the expectation of  $Z$  under  $D_{\text{KL}}$  exponentially converges to the expectation of  $Z^*$ ,  
 498 i.e.,  $\gamma$ -contraction.

499

□

500 We show that  $\mathcal{W}_{c,\varepsilon}(\mu, \nu)$  is NOT scale sensitive. Firstly, we denote  $\Pi^2$  as the optimal joint distribution  
 501 for  $(U, V)$  and thus we write the explicit form of Sinkhorn divergence  $W_{c,\varepsilon}(U, V)$  between two  
 502 random variables  $U$  and  $V$ :

503 \*\*\*

$$W_{c,\varepsilon}(U, V) = \text{KL}(\Pi^2 | \mathcal{K}) \quad (32)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \Pi^2(x, y) \log \frac{\Pi^2(x, y)}{\frac{1}{Z_2} e^{-\frac{c(x,y)}{\varepsilon}} \mu(x) \nu(y)} dx dy, \quad (33)$$

504 \*\*\*

505 where the normalization factor  $Z_2$  for the Gibbs kernel  $\mathcal{K}$  is  $Z_2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{c(x,y)}{\varepsilon}} \mu(x) \nu(y) dx dy$   
 506 and  $\mu(x), \nu(y)$  are the marginal density function of  $U$  and  $V$  with respect to  $x$  and  $y$ . We  
 507 also denote  $\Pi^1$  as the optimal joint distribution for  $(aU, aV)$ . A key proof element is about  
 508 the Gibbs kernel  $\mathcal{K}$ . By definition, the pdf of  $\mathcal{K}(U, V) \propto e^{-\frac{c(x,y)}{\varepsilon}} \mu(x) \nu(y)$ . After a scal-  
 509 ing transformation, the pdf of  $aU$  and  $aV$  with respect to  $x$  and  $y$  would be  $\frac{1}{a} \mu(\frac{x}{a})$  and  
 510  $\frac{1}{a} \nu(\frac{y}{a})$ . Thus  $\mathcal{K}(2U, 2V) \propto e^{-\frac{c(x,y)}{\varepsilon}} \frac{1}{a} \mu(\frac{x}{a}) \frac{1}{a} \nu(\frac{y}{a})$ . The new normalization factor  $Z_1$  is  $Z_1 =$   
 511  $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{a^2} e^{-\frac{c(x',y')}{\varepsilon}} \mu(x'/a) \nu(y'/a) dx' dy' = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{c(ax,ay)}{\varepsilon}} \mu(x) \nu(y) dx dy$ , the cost func-  
 512 tion of which is different from  $Z_2$ . For  $\Pi^2(U, V)$ , the scaled pdf of  $\Pi^2(aU, aV)$  would be  $\frac{1}{a^2} \Pi^2(\frac{x}{a}, \frac{y}{a})$ .  
 513 Then we have the following results:

514 \*\*\*

$$\mathcal{W}_{c,\varepsilon}(aU, aV) = \text{KL}(\Pi^1 | \mathcal{K}) \quad (34)$$

$$\leq \text{KL}(\Pi^2 | \mathcal{K}) \quad (35)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{1}{a^2} \Pi^2\left(\frac{x'}{a}, \frac{y'}{a}\right) \log \frac{\frac{1}{a^2} \Pi^2\left(\frac{x'}{a}, \frac{y'}{a}\right)}{\frac{1}{a^2} \frac{1}{Z_1} e^{-\frac{c(x',y')}{\varepsilon}} \mu\left(\frac{x'}{a}\right) \nu\left(\frac{y'}{a}\right)} dx' dy', \quad (36)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \Pi^2(x, y) \log \frac{\Pi^2(x, y)}{\frac{1}{Z_1} e^{-\frac{c(ax,ay)}{\varepsilon}} \mu(x) \nu(y) \frac{Z_2}{Z_1}} dx dy, \quad (37)$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \Pi^2(x, y) \left( \log \frac{\Pi^2(x, y)}{\frac{1}{Z_1} e^{-\frac{c(ax,ay)}{\varepsilon}} \mu(x) \nu(y)} + \log \frac{Z_1}{Z_2} \right) dx dy, \quad (38)$$

$$\stackrel{c=-k_\alpha, a \leq 1}{\leq} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \Pi^2(x, y) \log \frac{\Pi^2(x, y)}{\frac{1}{Z_1} e^{-\frac{c(x,y)}{\varepsilon}} \mu(x) \nu(y)} dx dy + \log \frac{Z_1}{Z_2} \cdot 1, \quad (39)$$

$$= \mathcal{W}_{c,\varepsilon}(U, V) + \Delta_{\mu,\nu}^c(a), \quad (40)$$

515 \*\*\*

516 where the second positive term  $\Delta_{\mu,\nu}^c(a) = \log \frac{Z_1}{Z_2}$  satisfies  $\Delta_{\mu,\nu}^c(a) \rightarrow 0$  as  $a \rightarrow 1$  (in practice  $\gamma$   
 517 is very close to 1). The second inequality holds for the general  $\varepsilon$  because for the unrectified kernel  
 518  $k_\alpha = -\|x - y\|^\alpha$  **with**  $a \leq 1$ , for any  $\varepsilon$  and  $x, y$  we have

$$(ax - ay)^\alpha \leq |a|^\alpha (x - y)^\alpha \leq (x - y)^\alpha e^{-\frac{c(ax,ay)}{\varepsilon}} \geq e^{-\frac{c(x,y)}{\varepsilon}} \Rightarrow e^{-\frac{c(ax,ay)}{\varepsilon}} \mu(x) \nu(y) \geq e^{-\frac{c(x,y)}{\varepsilon}} \mu(x) \nu(y)$$

519 However, under this condition,  $Z_1 \geq Z_2$  and thus  $\Delta_{\mu,\nu}^c(a) \geq 0$ , but  $\Delta_{\mu,\nu}^c(a) \rightarrow 0$  as  $a \rightarrow 1$  (in  
520 practice  $\gamma$  is very close to 1). We think there is indeed a gap between a (close) non-expansion property  
521 of Sinkhorn divergence and the empirical success of SinkhornDRL algorithm. The inequality is  
522 established based on the unrectified kernel, but it is tricky to find the contrative property for Sinkhorn  
523 divergence with the Gaussian kernel for any  $\varepsilon$  and  $x, y$ . Thus, it is fair that some counterexamples  
524 may exist for the non-contractive  $\mathfrak{T}^\pi$  under Sinkhorn divergence, which is also consistent with the  
525 counterexample MMD with Gaussian kernel (when  $\varepsilon \rightarrow \infty$ ).

526 Now we show that  $\mathcal{W}_{c,\varepsilon}$  is shift invariant:

$$\begin{aligned} \mathcal{W}_{c,\varepsilon}(A + X, A + Y) &= \int_{-\infty}^{\infty} \Pi^*(x - A, y - A) \log \frac{\Pi^*(x - A, y - A)}{\frac{1}{Z} e^{-\frac{c(x-A, y-A)}{\varepsilon}}} dx dy \\ &= \mathcal{W}_{c,\varepsilon}(X, Y). \end{aligned} \quad (41)$$

527 According to the equation of  $\overline{\mathcal{W}}_{c,\varepsilon}$ , it holds the same properties as  $\mathcal{W}_{c,\varepsilon}$ , i.e., shift invariant and scale  
528 sensitive. Thus, we derive the convergence of distributional Bellman operator  $\mathfrak{T}^\pi$  under the supreme  
529 form of  $\overline{\mathcal{W}}_{c,\varepsilon}$ , i.e.,  $\overline{\mathcal{W}}_{c,\varepsilon}^\infty$ :

$$\begin{aligned} &\overline{\mathcal{W}}_{c,\varepsilon}^\infty(\mathfrak{T}^\pi Z_1, \mathfrak{T}^\pi Z_2) \\ &= \sup_{s,a} \overline{\mathcal{W}}_{c,\varepsilon}(\mathfrak{T}^\pi Z_1(s, a), \mathfrak{T}^\pi Z_2(s, a)) \\ &= \overline{\mathcal{W}}_{c,\varepsilon}(R(s, a) + \gamma Z_1(s', a'), R(s, a) + \gamma Z_2(s', a')) \\ &\leq \overline{\mathcal{W}}_{c,\varepsilon}(Z_1(s', a'), Z_2(s', a')) + \Delta_{s',a',s,a}^{-k_\alpha}(\gamma) \\ &\leq \sup_{s',a'} \overline{\mathcal{W}}_{-k_\alpha,\varepsilon}(Z_1(s', a'), Z_2(s', a')) + \sup_{s,a,s',a'} \Delta_{s',a',s,a}^{-k_\alpha}(\gamma) \\ &= \overline{\mathcal{W}}_{-k_\alpha,\varepsilon}^\infty(Z_1, Z_2) + \Delta(\gamma) \end{aligned} \quad (42)$$

530 where the first inequality comes from the scale sensitivity proof, and we denote  
531  $\sup_{s,a,s',a'} \Delta_{s',a',s,a}^{-k_\alpha}(\gamma) = \Delta(\gamma)$  for short. Since  $\Delta(\gamma) \rightarrow 0$  as  $\gamma \rightarrow 1$ , we can conclude that  
532  $\mathfrak{T}^\pi$  is **closely** a non-expansive operator regardless of the cost function form  $c$  when  $\varepsilon \in (0, \infty)$ . The  
533  $\gamma$ -contraction of the expectation of  $Z^\pi$  can be similarly proved as the KL divergence in Lemma 1.  $\square$

## 534 C Proof of Proposition 1 and Corollary 1

535 *Proof.* As we leverage  $\Pi^*$  to denote the optimal  $\Pi$  by evaluating the Sinkhorn divergence via  
536  $\min_{\Pi \in \Pi(\mu, \nu)} \overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu; k)$ , the Sinkhorn divergence can be composed in the following form:

$$\begin{aligned} &\overline{\mathcal{W}}_{c,\varepsilon}(\mu, \nu; k) \\ &= 2\text{KL}(\Pi^*(\mu, \nu) | \mathcal{K}_{-k}(\mu, \nu)) - \text{KL}(\Pi^*(\mu, \mu) | \mathcal{K}_{-k}(\mu, \mu)) - \text{KL}(\Pi^*(\nu, \nu) | \mathcal{K}_{-k}(\nu, \nu)) \\ &= 2(\mathbb{E}_{X,Y} [\log \Pi^*(\mu, \nu)]) + \frac{1}{\varepsilon} \mathbb{E}_{X,X'} [c(X, Y)] - (\mathbb{E}_{X,X'} [\log \Pi^*(\mu, \mu)] + \frac{1}{\varepsilon} \mathbb{E}_{X,Y} [c(X, Y)]) \\ &\quad - (\mathbb{E}_{Y,Y'} [\log \Pi^*(\nu, \nu)] + \frac{1}{\varepsilon} \mathbb{E}_{Y,Y'} [c(Y, Y')]) \\ &= \mathbb{E}_{X,X',Y,Y'} \left[ \log \frac{(\Pi^*(X, Y))^2}{\Pi^*(X, X')\Pi^*(Y, Y')} \right] + \frac{1}{\varepsilon} (\mathbb{E}_{X,X'} [k(X, X')] + \mathbb{E}_{Y,Y'} [k(Y, Y')] - 2\mathbb{E}_{X,X'} [k(X, Y)]) \\ &= \mathbb{E}_{X,X',Y,Y'} \left[ \log \frac{(\Pi^*(X, Y))^2}{\Pi^*(X, X')\Pi^*(Y, Y')} \right] + \frac{1}{\varepsilon} \text{MMD}_{-c}^2(\mu, \nu) \end{aligned} \quad (43)$$

537 where the cost function  $c$  in the Gibbs distribution  $\mathcal{K}$  is minus Gaussian kernel, i.e.,  $c(x, y) =$   
538  $-k(x, y) = e^{-(x-y)/(2\sigma^2)}$ . Till now, we have shown the result in Corollary 1.

539 Next, we use Taylor expansion to prove the moment matching of MMD. Firstly, we have the following  
 540 equation:

$$\begin{aligned} \text{MMD}_{-c}^2(\mu, \nu) &= \mathbb{E}_{X, X'} [k(X, X')] + \mathbb{E}_{Y, Y'} [k(Y, Y')] - 2\mathbb{E}_{X, X'} [k(X, Y)] \\ &= \mathbb{E}_{X, X'} [\phi(X)^\top \phi(X')] + \mathbb{E}_{Y, Y'} [\phi(Y)^\top \phi(Y')] - 2\mathbb{E}_{X, X'} [\phi(X)^\top \phi(Y)] \quad (44) \\ &= \mathbb{E} \|\phi(X) - \phi(Y)\|^2 \end{aligned}$$

541 We expand the Gaussian kernel via Taylor expansion, i.e.,

$$\begin{aligned} k(x, y) &= e^{-(x-y)^2/(2\sigma^2)} \\ &= e^{-\frac{x^2}{2\sigma^2}} e^{-\frac{y^2}{2\sigma^2}} e^{\frac{xy}{\sigma^2}} \\ &= e^{-\frac{x^2}{2\sigma^2}} e^{-\frac{y^2}{2\sigma^2}} \sum_{n=0}^{\infty} \frac{1}{\sqrt{n!}} \left(\frac{x}{\sigma}\right)^n \frac{1}{\sqrt{n!}} \left(\frac{y}{\sigma}\right)^n \quad (45) \\ &= \sum_{n=0}^{\infty} e^{-\frac{x^2}{2\sigma^2}} \frac{1}{\sqrt{n!}} \left(\frac{x}{\sigma}\right)^n e^{-\frac{y^2}{2\sigma^2}} \frac{1}{\sqrt{n!}} \left(\frac{y}{\sigma}\right)^n \\ &= \phi(x)^\top \phi(y) \end{aligned}$$

542 Therefore, we have

$$\begin{aligned} \text{MMD}_{-c}^2(\mu, \nu) &= \sum_{n=0}^{\infty} \frac{1}{\sigma^{2n} n!} \left( \mathbb{E}_{x \sim \mu} \left[ e^{-x^2/(2\sigma^2)} x^n \right] - \mathbb{E}_{x \sim \nu} \left[ e^{-x^2/(2\sigma^2)} x^n \right] \right)^2 \quad (46) \\ &= \sum_{n=0}^{\infty} \frac{1}{\sigma^{2n} n!} \left( \tilde{M}_n(\mu) - \tilde{M}_n(\nu) \right)^2 \end{aligned}$$

543  $\tilde{M}_n(\mu) = \mathbb{E}_{x \sim \mu} \left[ e^{-x^2/(2\sigma^2)} x^n \right]$ , and similarly for  $\tilde{M}_n(\nu)$ . The conclusion is the same as the  
 544 moment matching in [17]. Finally, due to the equivalence of  $\bar{\mathcal{W}}_{c,\varepsilon}(\mu, \nu)$  after multiplying  $\varepsilon$ , we have

$$\begin{aligned} \bar{\mathcal{W}}_{c,\varepsilon}(\mu, \nu; k) &:= \text{MMD}_{-c}^2(\mu, \nu) + \varepsilon \mathbb{E} \left[ \frac{(\Pi^*(X, Y))^2}{\Pi^*(X, X') \Pi^*(Y, Y')} \right] \quad (47) \\ &= \sum_{n=0}^{\infty} \frac{1}{\sigma^{2n} n!} \left( \tilde{M}_n(\mu) - \tilde{M}_n(\nu) \right)^2 + \varepsilon \mathbb{E} \left[ \frac{(\Pi^*(X, Y))^2}{\Pi^*(X, X') \Pi^*(Y, Y')} \right], \end{aligned}$$

545 This result is also equivalent to Theorem 1, where  $\Pi^*$  would degenerate to  $\mu \otimes \nu$  as  $\varepsilon \rightarrow +\infty$ . In  
 546 that case, the first regularization term would vanish, and thus the Sinkhorn divergence degrades to a  
 547 MMD loss, i.e.,  $\text{MMD}_{-c}^2(\mu, \nu)$ .

548 □

## 549 D Human-normalized Scores

550 Our implementation is based on [25] and all the experimental settings, including parameters are  
 551 identical to the distributional RL baselines implemented by [25]. The main results about mean and

	Mean	Median	>Human	>DQN
DQN	173 %	49 %	17	0
C51	309 %	77 %	26	42
QR-DQN-1	430 %	104 %	31	47
MMDQN	<b>600 %</b>	94 %	27	43
SinkhornDRL	<u>570 %</u>	<u>89 %</u>	27	42

Table 2: Mean and median of best human-normalized scores across 55 Atari 2600 games. The results for all considered algorithms are averaged over 3 seeds.



552 median human-normalized scores of all considered distributional RL algorithms are reported in  
 553 Table 2. Note that our implementation is based on Pytorch, and thus the results in Table 2 are not  
 554 exactly same as results implemented based on Dopamine framework [4]. However, Table 2 also  
 555 suggests that our SinkhornDRL algorithm can achieve almost state-of-the-art performance in terms  
 556 of mean human-normalized scores. We argue that although it seems that SinkhornDRL is on par with  
 557 MMD across all games, our algorithm significantly outperforms MMDDRL on a large amount of Atari  
 558 games, as suggested in Figure 2. The detailed comparison based on learning curves is also exhibited  
 559 in Appendix E.

## 560 E More experimental Results

561 We provide learning curves of DQN, QRDQN, C51, MMD and SinkhornDRL algorithms on all  
 562 55 Atari games in Figures 4 5 6 7 8 9. It illustrates that SinkhornDRL dramatically surpasses the  
 563 other distributional RL algorithms on a large amount of environments, e.g., Venture, Atlantis, Tennis  
 564 and SpaceInvader, and presents competitive performance or is only slightly inferior as opposed to  
 565 the state-of-the-art baselines on other games. Note that the average improvement of SinkhornDRL  
 566 on Venture game is significant owing to one to two times convergence of SinkhornDRL algorithm  
 567 over 3 seeds, while the other baselines do not converge over the considered seeds. Although this  
 568 improvement may also suffer from the instability issue, its occasional success for our SinkhornDRL  
 569 algorithm also presents huge potential on some complicated environments. We leave the further  
 570 exploration on the advantage and potential of SinkhornDRL algorithm as the future work.

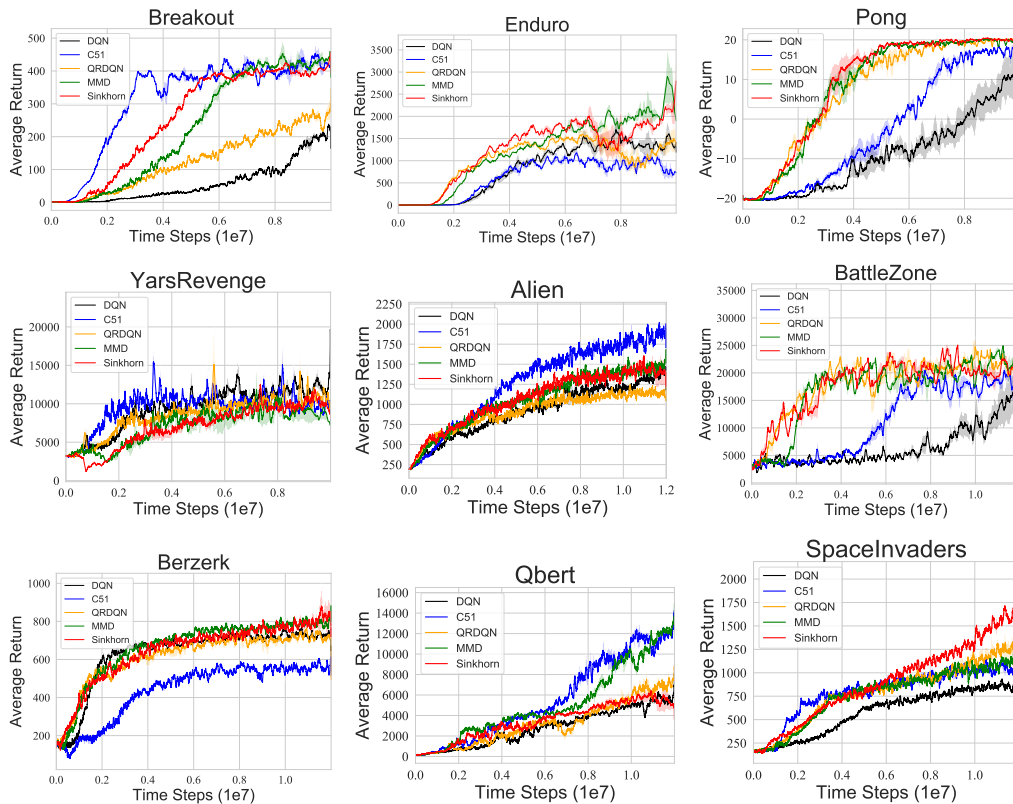


Figure 4: Performance of SinkhornDRL compared with DQN, C51, QRDQN and MMD on Breakout, Enduro, Pong, YarRevenge, Alien, BattleZone, Berzerk, Qbert and SpaceInvader.

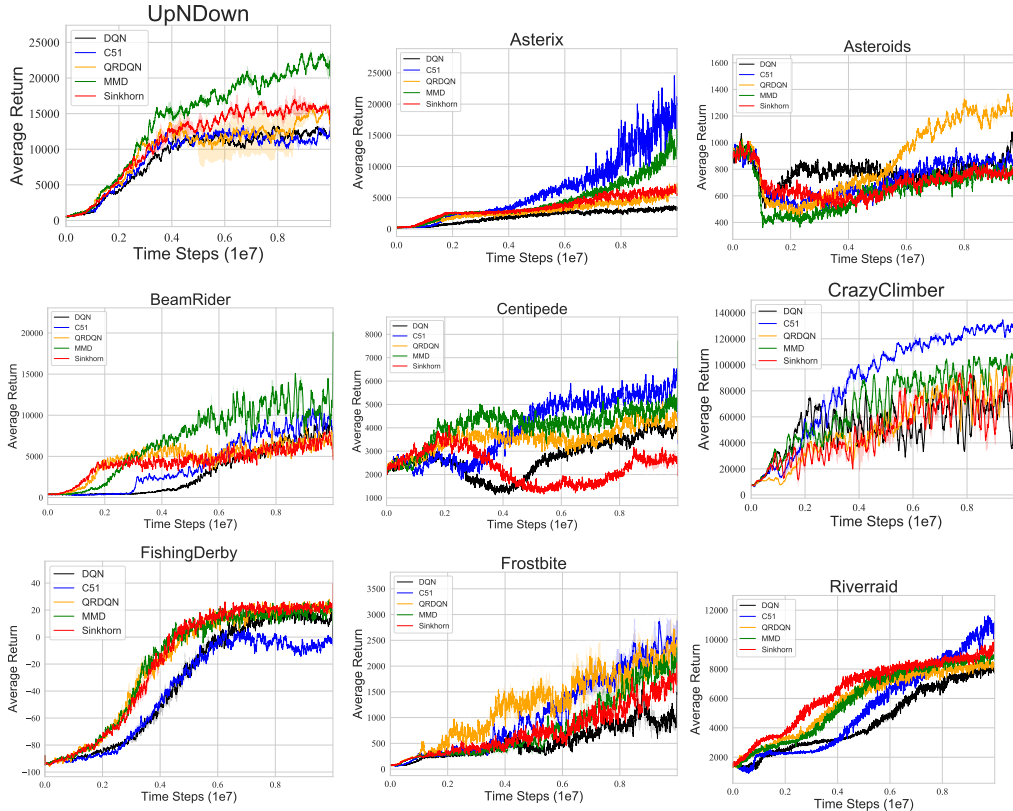


Figure 5: Performance of SinkhornDRL compared with DQN, C51, QRDQN and MMD on UpN-Down, Asterix, Asteroids, BeamRider, Centipede, FishingDerby, Frostbite and Riverraid.

## 571 F Sensitivity Analysis and Computational Cost

### 572 F.1 More results in Sensitivity Analysis

573 From Figure 10 (a), we can observe that if we gradually decline  $\varepsilon$  to 0, SinkhornDRL's performance  
 574 tends to QR-DQN. Note that an overly small  $\varepsilon$  will lead to a trivial almost 0  $\mathcal{K}_{i,j}$  in Sinkhorn iteration  
 575 in Algorithm 2, and will cause  $\frac{1}{0}$  numerical instability issue for  $a_l$  and  $b_l$  in Line 5 of Algorithm 2.  
 576 Due to this reason, the performance of SinkhornDRL with  $\varepsilon = 0.1$  or  $0.075$  declines as the training  
 577 proceeds, and eventually converges to the average return that QR-DQN achieves. In addition, we also  
 578 conducted experiments on Seaquest, the similar result is also observed in Figure 11. The performance  
 579 of SinkhornDRL is robust when  $\varepsilon = 10, 100, 500$  and a small  $\epsilon = 1$  tends to worsen the performance.

580 Moreover, for breakout, if we increase  $\varepsilon$ , the performance of SinkhornDRL tends to that of MMDDRL  
 581 as suggested in Figure 10 (b). It is also noted that an overly large  $\varepsilon$  will let the  $\mathcal{K}_{i,j}$  explode to  $\infty$ .  
 582 This also leads to numerical instability issue in Sinkhorn iteration in Algorithm 2.

583 In summary, the trend of SinkhornDRL to close MMDDRL and QR-DQN if we increase or decrease  $\varepsilon$ ,  
 584 respectively, provides strong empirical evidence to demonstrate the theoretical relationships between  
 585 Sinkhorn divergence and MMD / Wasserstein distance, although an overly large or small  $\varepsilon$  will lead  
 586 to numerical instability issue.

### 587 F.2 Comparison with the Computational Cost

588 We evaluate the computational time every 10,000 iterations across the whole training process of  
 589 all considered distributional RL algorithms and make a comparison in Figure 12. It suggests that  
 590 SinkhornDRL indeed increases around 50% computation cost compared with QR-DQN and C51,  
 591 but only slightly increases the the cost in contrast to MMDDRL on both Breakout and Qbert

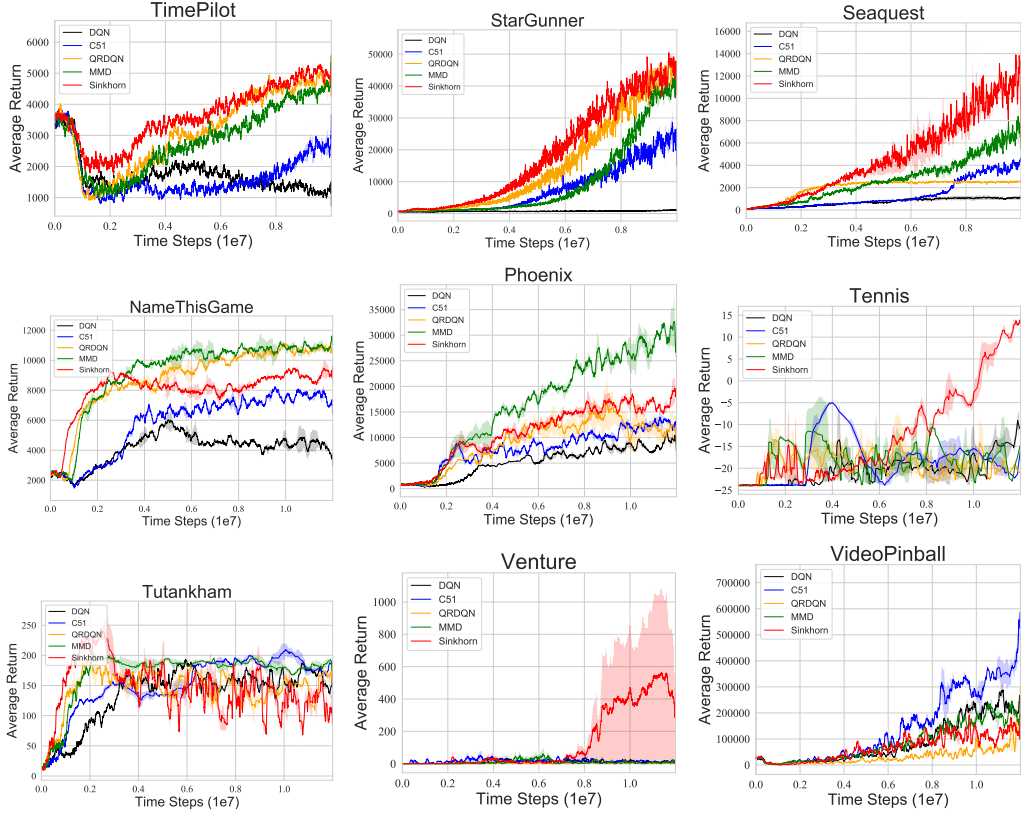


Figure 6: Performance of SinkhornDRL compared with DQN, C51, QRDQN and MMD on TimePilot, StarGunner, Seaquest, NameThisGame, Phoenix, Tennis, Tutankham, Venture and VideoPinball.

592 games. We argue that this additional computational burden can be tolerant in view of the significant  
 593 outperformance of SinkhornDRL in a large amount of environments.

594 In addition, we also find that the number of Sinkhorn iterations  $L$  is negligible to the computation cost,  
 595 while an overly large samples  $N$ , e.g., 500, will lead to a large computational burden as illustrated in  
 596 Figure 13. This can be intuitively explained as the computation complexity of the cost function  $c_{i,j}$  is  
 597  $\mathcal{O}(N^2)$  in SinkhornDRL, which is particularly heavy in computation if  $N$  is large enough.

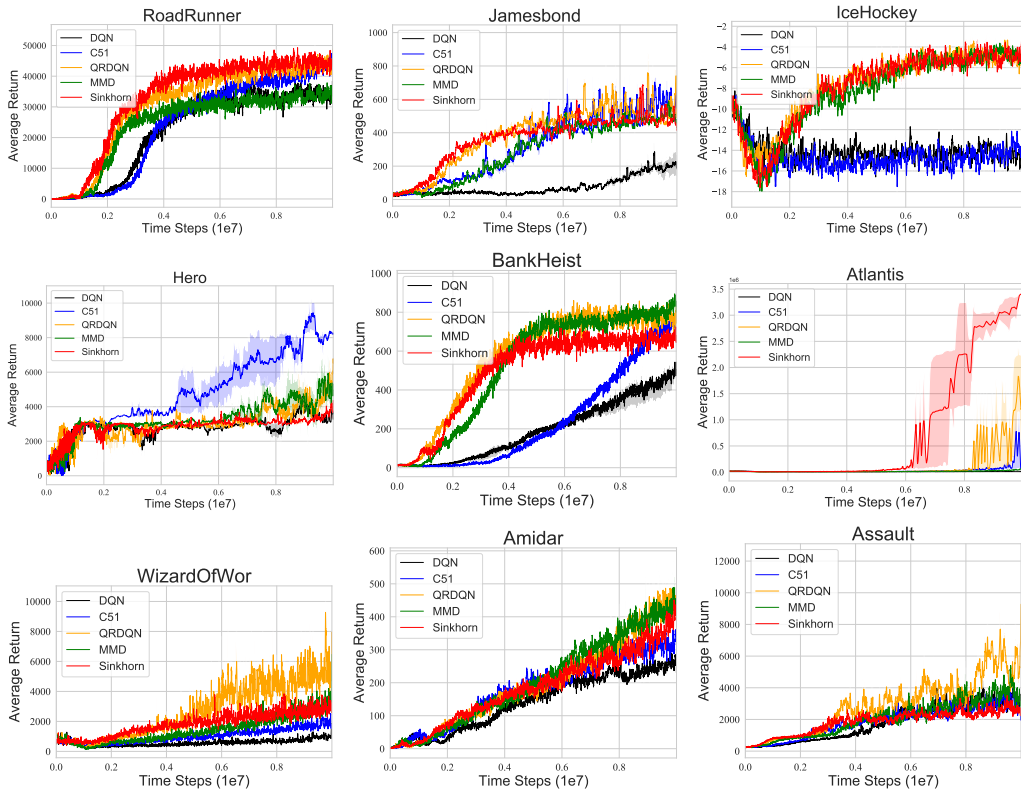


Figure 7: Performance of SinkhornDRL compared with DQN, C51, QRDQN and MMD on Road-Runner, Jamesbond, IceHockey, Hero, BankHeist, Atlantis, WizardOfWor, Amidar and Assault.

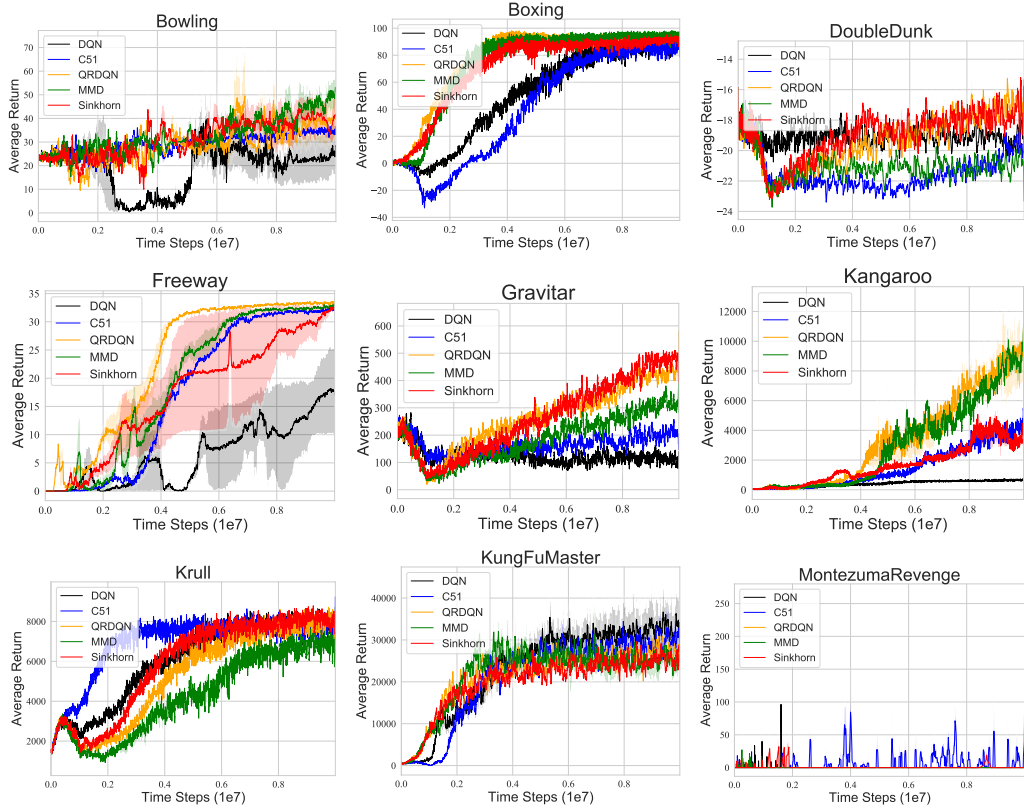


Figure 8: Performance of SinkhornDRL compared with DQN, C51, QRDQN and MMD on Bowling, Boxing, DoubleDunk, Freeway, Gravitar, Kangaroo, Krull, KunFuMaster and MontezumaRevenge.

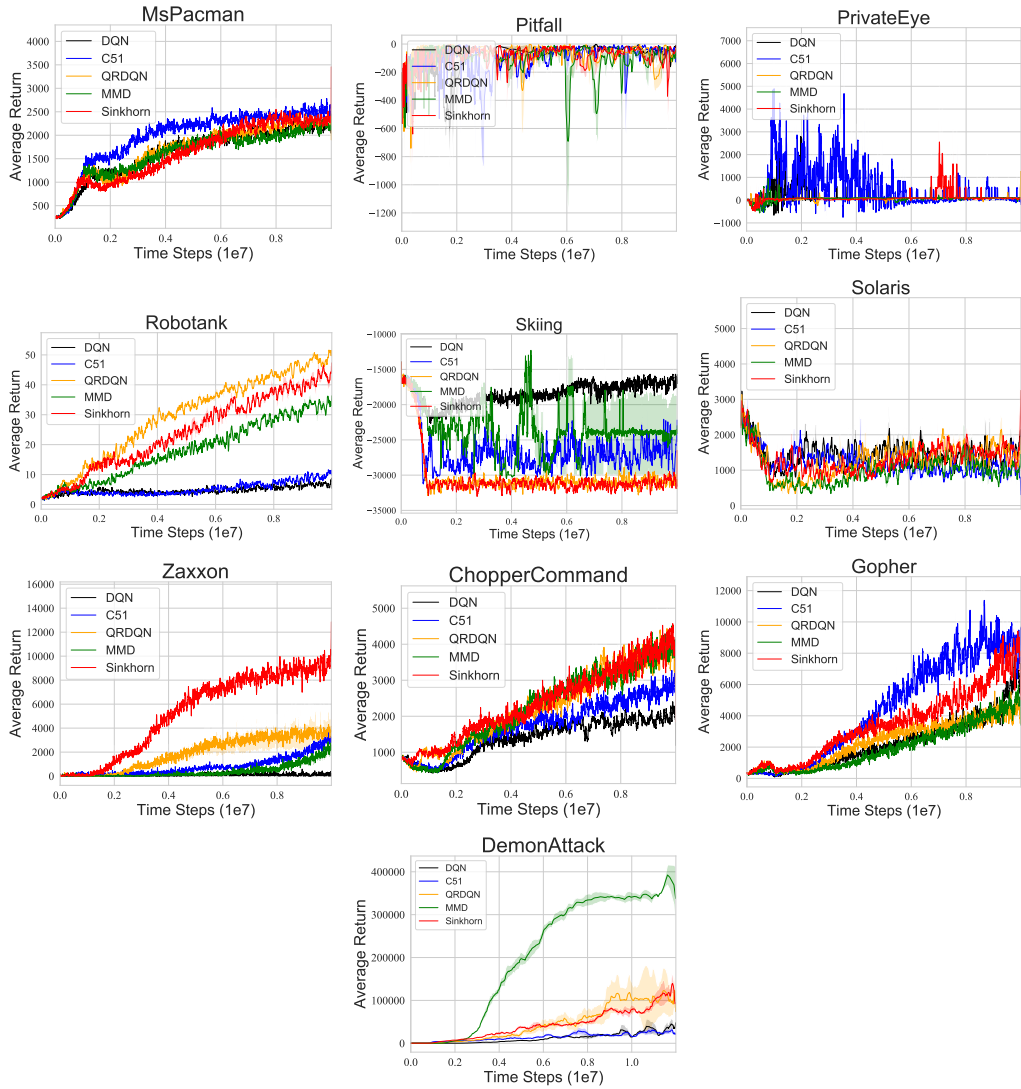
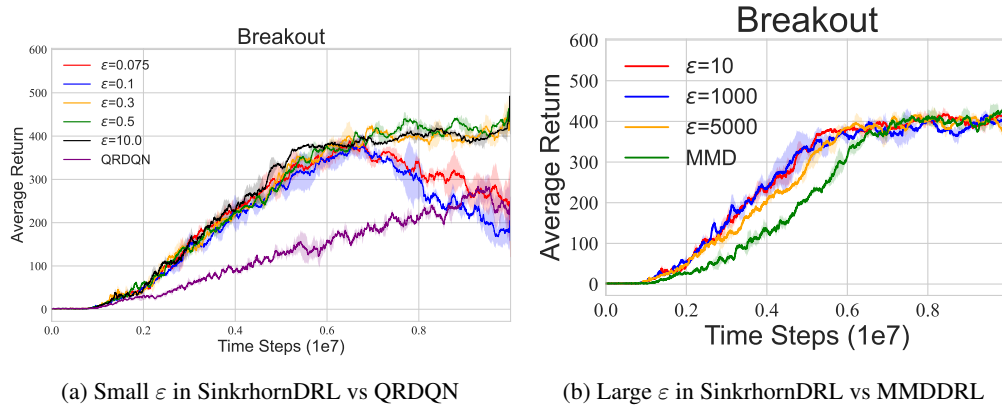


Figure 9: Performance of SinkhornDRL compared with DQN, C51, QRDQN and MMD on MsPacman, Pitfall, PrivateEye, Robotank, Skiing, Solaris, Zaxxon, ChopperCommand, Gopher and DemonAttack.



(a) Small  $\epsilon$  in SinkhornDRL vs QRDQN

(b) Large  $\epsilon$  in SinkhornDRL vs MMDDRL

Figure 10: (Left) Sensitivity analysis w.r.t. a small level of  $\epsilon$  SinkhornDRL to compare with QR-DQN that approximates Wasserstein distance on Breakout. (Right) Sensitivity analysis w.r.t. a large level of  $\epsilon$  SinkhornDRL algorithm to compare with MMDDRL on Breakout. All learning curves are reported over 2 seeds.

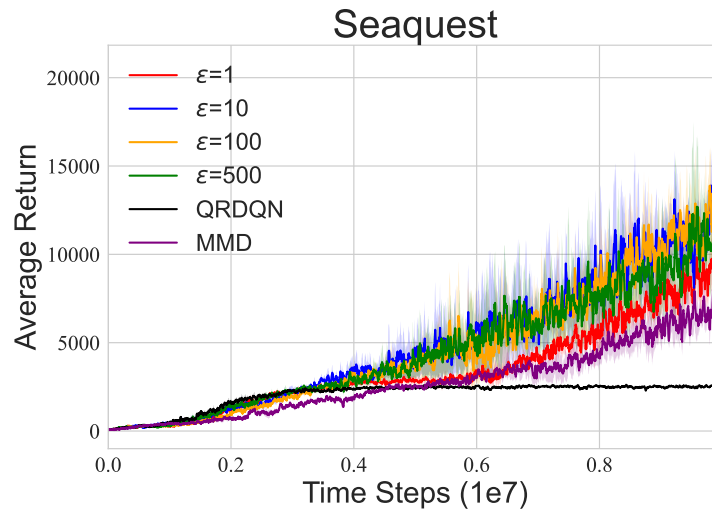


Figure 11: Sensitivity analysis w.r.t.  $\epsilon$  SinkhornDRL to compare with QR-DQN and MMD on Seaquest. All learning curves are reported over 3 seeds.

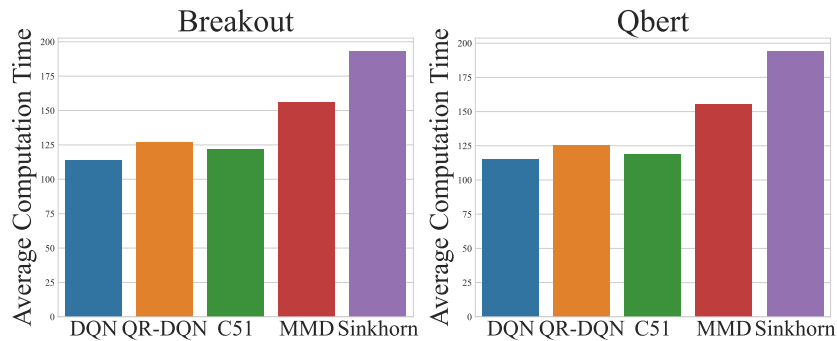


Figure 12: Average computational cost per 10,000 iterations of all considered distributional RL algorithm, where we select  $\epsilon = 10$ ,  $L = 10$  and number of samples  $N = 200$  in SinkhornDRL algorithm.

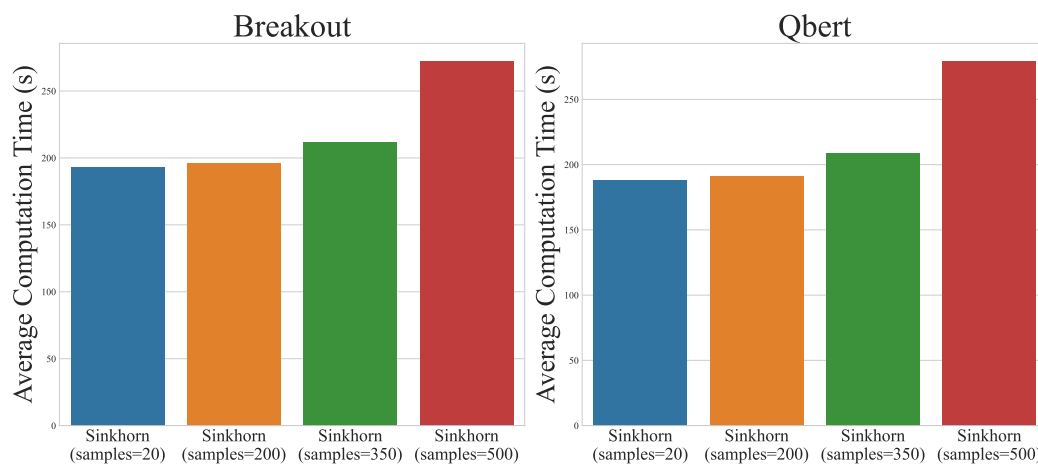


Figure 13: Average computational cost per 10,000 iterations of SinkhornDRL algorithm over different samples.