

1 Thank you for all the valuable feedback. One common comment is that some results may not be comparable to the
 2 original paper. Since we are not researchers from the industry and do not have strong computing power, we set the
 3 number of quantiles to be 100 and evaluate all methods using 40 million training frames for the experiment. We
 4 will definitely try to add the **complete comparison (200M, 57 games)** and more evaluation results (**stochastic Atari,**
 5 **different quantile numbers etc**) into the final version of the paper. We are also working on a modified version of the
 6 NC architecture which can be extended to IQN, and its comparison with the baseline IQN (using Google 'Dopamine')
 7 are provided in Figure (a) (last two with exploration). Our response to other specific comments are provided as follows.

8 **Reviewer #1:**

9 **Q1. Can authors explain why only 49 out of 57 games are used for evaluation?**

10 We choose the 49 Atari game (initially proposed by Mnih et al., 2015) since our main contrast DLTV is evaluated using
 11 the same environment. But we will definitely include all the 57 games in the final version.

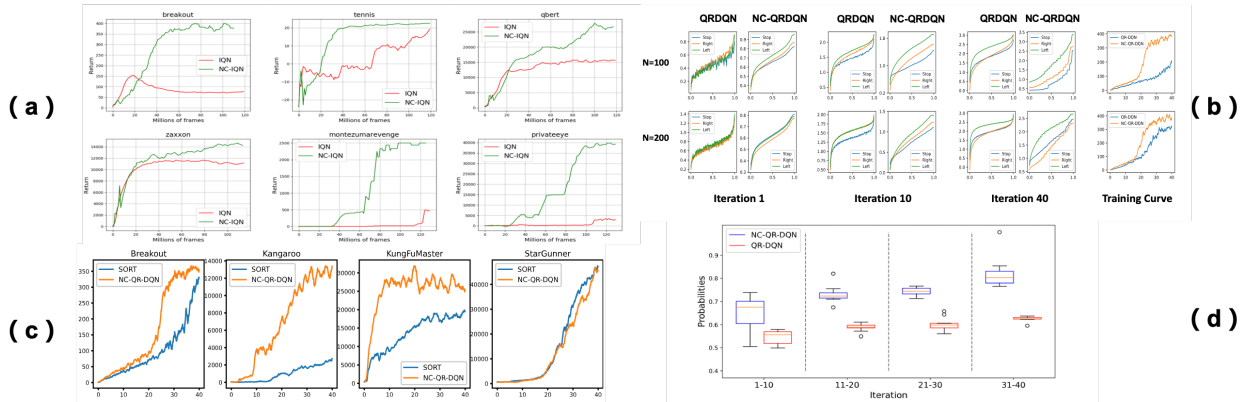
12 **Q2. If the authors can commit to providing detailed evaluation results over different number of quantiles.**

13 We will provide detailed evaluations in the final version of the paper. Figure (b) compares $N = 100, 200$ on Breakout.
 14 The crossing issue is more severe with smaller N at early stage, and the improvement of our approach is more significant.

15 **Q3. Hard exploration games: I would love to see complete results on full 57 games.**

16 Sorry for the misunderstanding. Figure 5(a) in the paper is actually the comparison between our method (with
 17 exploration) and DLTV on all the 49 games. The improvement in Motezuma's Revenge, and Private Eye are 6545%
 18 (60M) and 25% (40M) respectively. We will improve the presentation of this part in the final version.

19 **Reviewer #2:**



20 **Q1. It seems much natural to simply sort the atoms before calculating the quantile Huber loss.**

21 As Figure (c) shows, our approach performs better than 'sort' with a higher optimal return. The key reason is that,
 22 by sorting, the same dimension of the network output may be paired to a different quantile location ' τ ' each time
 23 without employing the non-decreasing constraint, which highly decreases the training stability and efficiency. With the
 24 monotonicity restriction, we can make better use of global information in training.

25 **Q2. The idea of using cumulated softmax is very similar to the Fully Parameterized Quantile Function (FPQF).**

26 Sorry for ignoring this important reference. We will cite it in the final draft of the paper.

27 **Reviewer #3:**

28 **Q1. Please expand upon how the non-crossing fix improves model interpretability.**

29 Sometimes, the upper quantiles instead of the mean (Q-value) are of specific interests when examining some risk-
 30 appetite policy. Crossing issue will lead to awkward interpretation due to the abnormal ranking of the quantile points.
 31 Also, the non-crossing fix can bring a more precise estimation of distribution variance when doing exploration.

32 **Q2. A plot showing how the ranking of Q-function changes on a variety of states would be more convincing.**

33 We randomly pick 4000 states, and compute the probabilities that QRDQN or NC-QRDQN chooses the same action
 34 with the optimal policy within each of the four training period. As the boxplots in Figure (d) shows, our method
 35 performs much more stable especially in the early stage with an overall higher consistency with the optimal policy.

36 **Reviewer #4:**

37 **Q1. "When the sample size goes to infinity, ..." and the statement in line 122 don't seem like clear results.**

38 When N is fixed and sample size goes to infinity, $\theta_i(s, a)$ converges to the real quantile function F_Z^{-1} at level $\hat{\tau}_i$ as
 39 defined in line 105. Thus, the monotonicity of quantile function guarantees the monotonic constraint in equation (9).
 40 The projection operator Π_{W_1} is to find a Z_q in non-crossing space \mathcal{Z}_Q to minimize its W_1 distance from Z , which is
 41 equivalent to finding a model parameterized by θ_i 's under the monotonicity constraint.

42 **Q2. Do we require \mathcal{Z}_Q to be a Banach space, to ensure the proper convergence of the operator to the fixed point.**

43 Yes, we need \mathcal{Z}_Q to be a complete space to ensure the proper convergence.