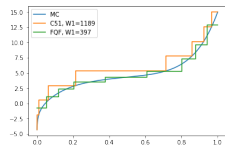
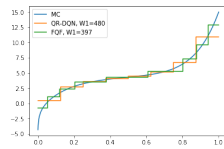


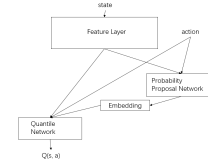
1 *General response:* Thanks a lot for your comments and suggestions! We will fix typos, add more details and re-organize  
2 the paper to improve clarity.



3 Figure 1



4 Figure 2



5 Figure 3

6 **(A) Toy Experiment for Distribution Approximation:** Figure 1 and 2 are toy cases of quantile function approximation  
7 using C51, QR-DQN and FQF (IQN should be viewed as sampling in a distribution instead of approximating it).  $W_1$   
8 loss suggests that FQF does have advantage in approximating distribution and thus achieve better estimation on values.

9 **(B) Experiment Details:** Figure 3 shows the general architecture of FQF. The learning rate of the probability network is  
10 set to 0.0001 but as long as its smaller than 0.001 FQF works quite well. We commit to add more experimental details in  
11 the camera-ready version. **(C) More Experiment under Different Settings:** We commit to compare FQF with IQN under  
12 different set of hyperparameters for ablative analysis in the camera-ready version. **(D) Inefficient Hyperparameter:** By  
13 ‘inefficient hyperparameter’ we mean the atom locations (uniformly distributed between -10 and 10) in C51, quantiles  
14 locations in QR-DQN and sampled quantiles locations in IQN (it still requires a distribution to sample from). Note that  
15 FQF requires only the number of quantiles. **(E) State-of-Art Performance:** Thanks for pointing out [Kapturowski et al.,  
16 2019]! Yes, FQF is the best among single-actor algorithms, i.e., non-distributed methods. We will clarify this in the  
17 new version and make our claim more accurate. Besides, we believe that combining advantages of distributional RL  
18 and distributed RL is an exciting direction for further research.

19 *To reviewer 1:* **(I) Experiment Details:** Please refer to **(B)** in general response. **(II) IQN v.s. FQF:** IQN does have the  
20 advantage that it can sample as many quantiles as possible, but with the same number of quantiles FQF gives much  
21 more accurate approximations. It is not difficult to extend FQF to support arbitrary number of optimal quantiles as  
22 in IQN: we can modify the quantile network into a recurrent network so that it generates a single  $\tau$  each time and  
23 takes state, action and previously outputted  $\tau$ s as input to generate next  $\tau$ . FQF finds quantiles that can most efficiently  
24 express a distribution and computes the expectation of Q with such expression. Thus, intuitively the quantiles generated  
25 in FQF should be better than sampled quantiles (as in IQN) in terms of quantile function approximation in most cases.  
26 We leave the theoretical analysis and comparison between IQN and FQF to future work.

27 *To reviewer 2:* **(I) More Experiment under Different Settings:** Please refer to **(C)** in general response. We will add  
28 experiments regarding different K and N in the new version. It is a great idea to investigate what impacts the quantiles  
29 generated in FQF, and we will include the related results in the appendix of the new version as well as the results  
30 using sticky actions. **(II) Different Weights v.s. Uniform Weights:** We actually emailed the authors of IQN regarding  
31 the weighting scheme and tested it. We all found that using different weights or uniform weights does not really  
32 impact performance. **(III) The Motivation for Minimizing  $L_w$ :** Why we minimize the square of the gradient instead of  
33 following the gradient, we agree that it should be explained. Although we cannot compute  $W_1$ , we derived its derivative  
34 for  $\tau$  in the appendix and directly following the gradient is the most intuitive way to minimize  $W_1$ . Unfortunately, we  
35 found following the gradient makes the probability proposal network very hard to train at the beginning when quantile  
36 outputs are unstable. It often causes  $\tau$ s to converge to something like (0,0,0,1,1,1,1), while using square (not very  
37 different from following the gradient actually) performs much better. This may be a result of the cumulative softmax  
38 architecture. We will check if we can figure out why and explain it in detail.

39 *To reviewer 3:* **(I) How well the distribution is approximated/Inefficient Hyper-parameters:** Please refer to **(A)/(D)** in  
40 the general response. **(II) Supporting Experiments:** We agree with you on adding supporting experiments other than  
41 performance, some of those experiments are shown in the general response section and we will add more to the new  
42 version. We also agree that performance measurements on variance is important for further distributional RL research  
43 and should be included. **(III) Why distribution matters?** The reason why distribution matters had been studied by  
44 several previous works: 1). An Analysis of Categorical Distributional Reinforcement Learning (Rowland et al. 2018)  
45 2). As Expected? An Analysis of Distributional Reinforcement Learning (Lyle and Bellemare, 2018). The general  
46 conclusion is that “the more complex the setting the less likely it is that distributional and expected RL algorithms  
47 will behave in the same way”. So it is necessary for distributional RL algorithms to approximate the full distribution,  
48 even when we only need its mean. As shown in our toy case, FQF does achieve better distribution approximation. **(IV)**  
49 “faster and more stable” We apologize for the misleading choice of word in “faster and more stable”. What we actually  
50 meant by ‘faster’ was higher sampling efficiency compared with IQN instead of training speed, and by stability we refer  
51 to smaller error bands. We will add numerical results to support this claim in the new version.