# 535514: Reinforcement Learning
# Lecture 24 — QR-DQN and IQN

Ping-Chun Hsieh

May 16, 2024

# Announcement

- No class next Monday (5/20) and next Thursday (5/23)

5月

| 日 | 一 | 二 | 三 | 四 | 五 | 六 |
|---|---|---|---|---|---|---|
| | | | 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| 26 | 27 | 28 | 29 | 30 | 31 | |

**Lec 24**

**No class** **No class**

**Lec 25 (SAC)** **Lec 26 (Inverse RL)**

6月

**Lec 27 (Model-based RL)** **Lec 28 (Offline MBRL)**

| 日 | 一 | 二 | 三 | 四 | 五 | 六 |
|---|---|---|---|---|---|---|
| | | | | | | 1 |
| 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 |
| 23 | 24 | 25 | 26 | 27 | 28 | 29 |
| 30 | | | | | | |

**Final presentation**

# On-Policy vs Off-Policy Methods

| | Policy Optimization | Value-Based | Model-Based | Imitation-Based |
|---|---|---|---|---|
| **On-Policy** | **Exact PG**<br>**REINFORCE (w/i baseline)**<br>**A2C**<br>**On-policy DAC**<br>**TRPO**<br>**Natural PG (NPG)**<br>**PPO-KL & PPO-Clip**<br>**RLHF by PPO-KL** | **Epsilon-Greedy MC**<br>**Sarsa**<br>**Expected Sarsa** | **Model-Predictive Control (MPC)**<br>**PETS** | **IRL**<br>**GAIL**<br>**IQ-Learn** |
| **Off-Policy** | **Off-policy DPG & DDPG**<br>**Twin Delayed DDPG (TD3)** | **Q-learning**<br>**Double Q-learning**<br>**DQN & DDQN**<br>**Rainbow**<br>**C51 / QR-DQN / IQN**<br>**Soft Actor-Critic (SAC)** | | |

# Quick Review: Distributional Bellman

▸ <u>Sample action-value $Z^\pi(s, a)$</u>: sample return if we start from state $s$ and take action $a$, and then follow policy $\pi$

$$Q^\pi(s, a) = \mathbb{E}[Z^\pi(s, a)] = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)\right]$$

▸ Distributional Bellman operator $B^\pi : \mathscr{Z} \to \mathscr{Z}$

$$B^\pi Z(s, a) \overset{D}{:=} r(s, a) + \gamma P^\pi Z(s, a)$$

where $P^\pi Z(s, a) \overset{D}{:=} Z(s', a')$

$$s' \sim P(\,\cdot \mid s, a), \ \ a' \sim \pi(\,\cdot \mid s')$$

▸ Distributional optimality operator $B*$: The $B^\pi$ resulting from a greedy policy $\pi$ (what does "greedy" mean here?)

# Quick Review: C51

$$a* = \arg\max_a \mathbb{E}[Z(s', a)]$$

(C1) <u>Categorical</u> distributions for parametrizing $Z_\theta(s, a)$



$V_{max}, V_{min}$

$Z(s', a^*)$

$p_1$
$p_0$
$p_9$

$z_0 z_1$
$z_9$

$r + \gamma Z$

Actions

Returns

$\psi$

C51

state

(C2) <u>Mimicking</u> $B*$ for learning with sample transitions $(s, a, r, s')$

(C3) <u>Cramer Projection</u> $\Phi$ for support mismatch caused by $B*Z_\theta(s, a)$

(C4) Minimize $L_{C51}(s, a, r, s'; \theta) := D_{KL}(\Phi B*Z_{\bar{\theta}}(s, a) \| Z_\theta(s, a))$

# QR-DQN

Dabney et al., Distributional Reinforcement Learning with Quantile Regression, AAAI 2018

# Quantile-Based Parametrization of $Z(s, a)$

▶ **Idea**: **Express** $Z(s, a)$ **using CDF (instead PDF)**

CDF $(P(Z(s, a) \leq z))$

$\tau_6 = 1$
$\tau_5$
$\tau_4$
$\tau_3$
$\tau_2$
$\tau_1$
$0$

$v_1$ $v_2$ $v_3$ $v_4$ $v_5$ $v_6$

PMF $(P(Z(s, a) = v_i))$

*Quantile distribution*
$\hat{Z}(s, a) \approx Z(s, a)$

quantile
distribution

$v_i = F_Z^{-1}(\tau_i)$

$v_1$ $v_2$ $v_3$ $v_4$ $v_5$ $v_6$

(inverse CDF)

Quantile function: $F_Z^{-1}(\tau) := \inf\{z : P(Z \leq z) \geq \tau\}$

$P(Z_{(s,a)} \leq v_1) = \tau_1$

7

# A Comparison of NN Architecture



Actions

Returns

$f$

$\psi$

**DQN**

$Z(s,a)$

Actions

$z_0$ $z_1$ $z_2$ $z_3$ $z_4$

Returns

$\psi$

**C51**

5 quantiles

$Z(s,a)$

Actions

Returns

$\psi$

**QR-DQN**

$\nu_1$ $\nu_2$ $\nu_3$ $\nu_4$ $\nu_5$

8

# QR-DQN: Another Popular Distributional DQN

▸ **Q1**: How to express $Z(s, a)$? ✓

(D1) **Quantile** distributions for $Z_\theta(s, a)$



DQN          QR-DQN

▸ **Q2**: How to update $Z(s, a)$ during training? ✓

(D2) Mimicking $\boxed{B*}$ for learning with sample transitions $(s, a, r, s')$

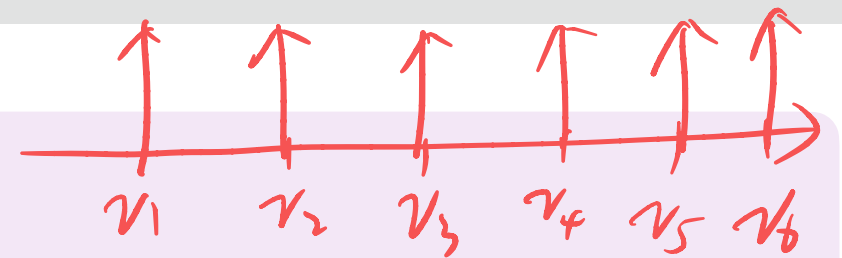(D3) Minimize $L_{QR}(s, a, r, s'; \theta) := D(B*Z_{\bar\theta}(s, a) \| Z_\theta(s, a))$

*No Cramer projection required!*

# Quantile Regression DQN (Formally)

Step 1: Initialize $Z_\theta(s, a)$ and initial state $s_0$

Step 2: For each step $t = 0, 1, 2, \cdots$

Select $a_t$ using $\varepsilon$-greedy w.r.t $Q(s_t, a) \equiv \mathbb{E}[Z_\theta(s_t, a)]$

Observe $(r_{t+1}, s_{t+1})$ and store $(s_t, a_t, r_{t+1}, s_{t+1})$ in the buffer

Draw a mini-batch of samples $B$ from the replay buffer

Update $\theta$ by minimizing QR loss as follows:

$$\theta \leftarrow \theta - \alpha \nabla_\theta \sum_{(s,a,r,s') \in B} L_{QRDQN}(s, a, r, s'; \theta)$$

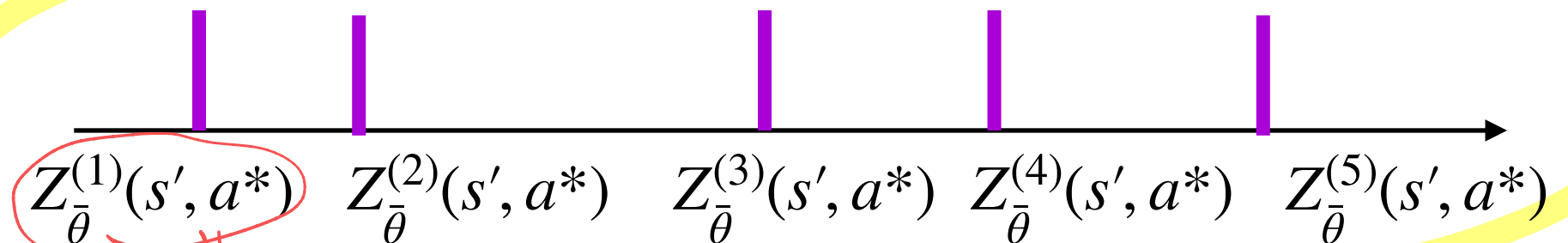Under quantile distributions, $\mathbb{E}[Z_\theta(s, a)] = \sum_{i=1}^{N} \frac{1}{N} Z_\theta^{(i)}(s, a)$

*(handwritten annotations: "quantile distribution" pointing to $Z_\theta$; arrows over a number line labeled $\nu_1, \nu_2, \nu_3, \nu_4, \nu_5, \nu_6$; red ellipse and underline around $\mathbb{E}[Z_\theta(s_t, a)]$)*

$$B^* Z(s,a) = r(s,a) + \gamma P^{\pi_{greedy}} \underline{Z}$$

▸ Here we presume a <u>greedy</u> policy w.r.t $Q$ function for $B*$ $\quad Z(s',a^*)$

▸ **Question**: Given only transitions $(s, a, r, s')$, how to enforce $B*$ to update $Z(s, a)$ on *quantile* distributions?

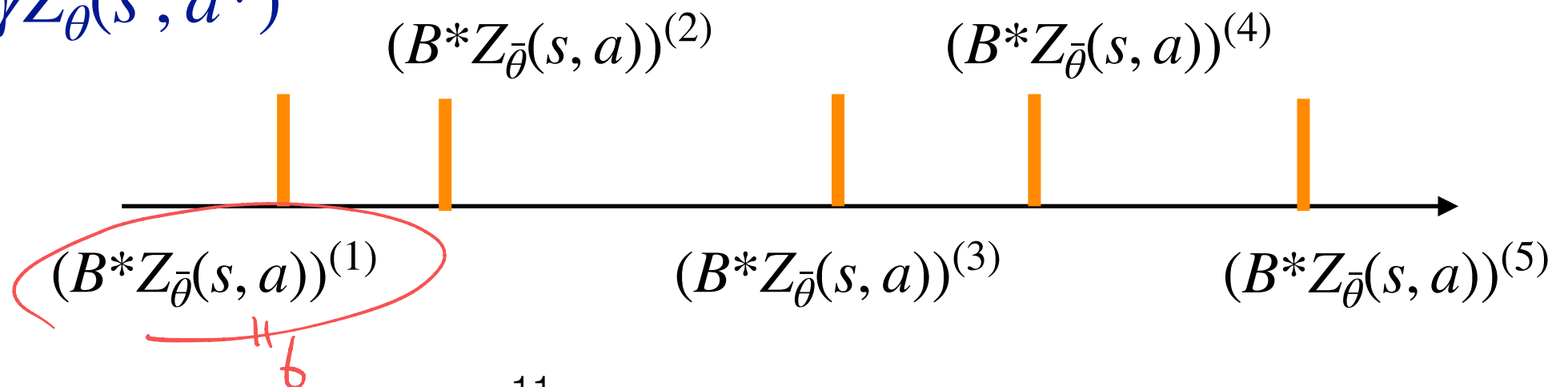$$a* = \arg\max_a Q(s', a) \equiv \arg\max_a \mathbb{E}[Z_{\bar{\theta}}(s', a)]$$

$Z_{\bar{\theta}}(s', a*)$

$Z_{\bar{\theta}}^{(1)}(s', a*) \quad\quad Z_{\bar{\theta}}^{(2)}(s', a*) \quad\quad Z_{\bar{\theta}}^{(3)}(s', a*) \quad Z_{\bar{\theta}}^{(4)}(s', a*) \quad\quad Z_{\bar{\theta}}^{(5)}(s', a*)$

"5

$$B^* Z_{\bar{\theta}}(s, a) = r + \gamma Z_{\bar{\theta}}(s', a*)$$

Suppose $\gamma = 0.9$

$\gamma = 1.5$

$(B^* Z_{\bar{\theta}}(s,a))^{(2)} \quad\quad\quad (B^* Z_{\bar{\theta}}(s,a))^{(4)}$

$(B^* Z_{\bar{\theta}}(s,a))^{(1)} \quad\quad\quad (B^* Z_{\bar{\theta}}(s,a))^{(3)} \quad\quad\quad (B^* Z_{\bar{\theta}}(s,a))^{(5)}$

"6

11

# (D3) Loss Function

- We still need to choose a "**dissimilarity**" function $D(\ \cdot\ ||\ \cdot\ )$ in
$$L_{QRDQN}(s, a, r, s'; \theta) := D(B^*Z_{\bar{\theta}}(s, a) || Z_{\theta}(s, a))$$

- There are many possibilities, e.g., total variation or KL divergence

- QR-DQN uses the quantile regression loss
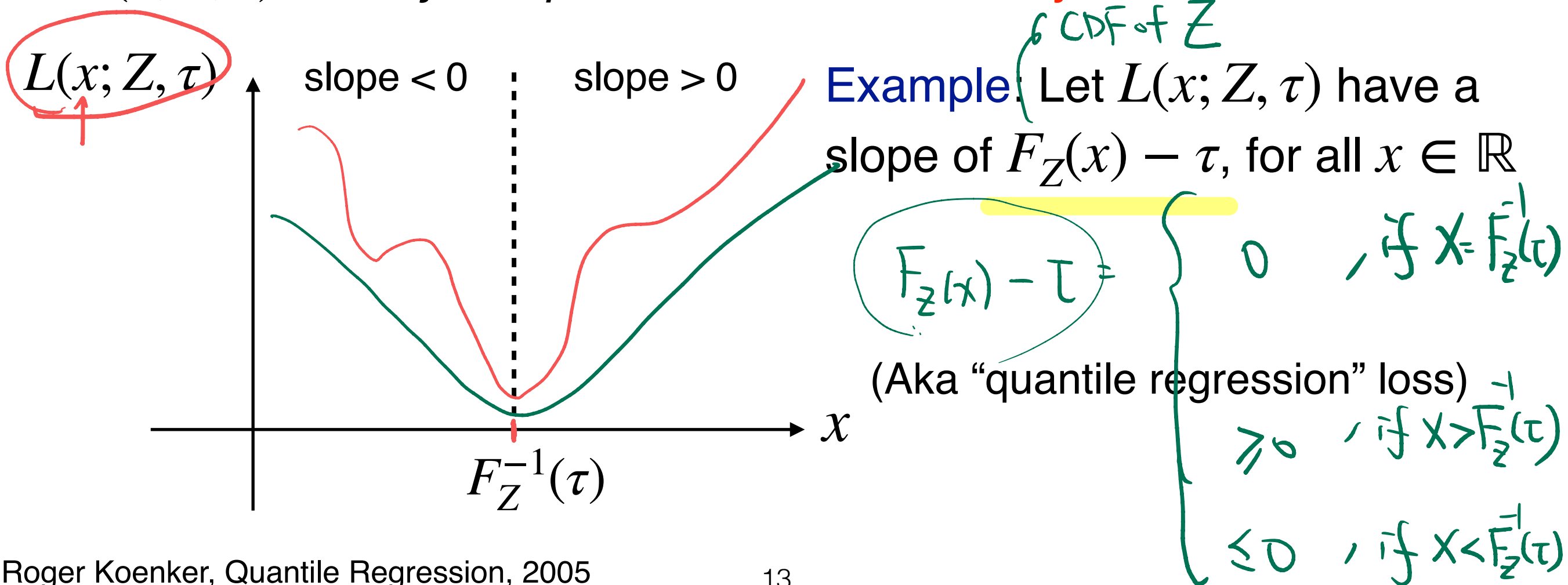  - Motivation: Both $B^*Z_{\bar{\theta}}(s, a)$ and $Z_{\theta}(s, a)$ are quantile distributions

# Quantile Regression Loss

· Given a random variable $Z$

· Goal: Find a quantile $F_Z^{-1}(\tau)$

▸ **Idea**: Finding a quantile $F_Z^{-1}(\tau)$ by minimizing loss $L(x; Z, \tau)$

$$F_Z^{-1}(\tau) = \arg\min_{x \in \mathbb{R}} L(x; Z, \tau)$$

▸ $L(x; Z, \tau)$ is *easy-to-optimize* when it is strictly convex



slope < 0      slope > 0

$L(x; Z, \tau)$

$F_Z^{-1}(\tau)$

$x$

↳ CDF of $Z$

**Example**: Let $L(x; Z, \tau)$ have a slope of $F_Z(x) - \tau$, for all $x \in \mathbb{R}$

$$F_Z(x) - \tau = \begin{cases} 0 & \text{, if } x = F_Z^{-1}(\tau) \\ > 0 & \text{, if } x > F_Z^{-1}(\tau) \\ \leq 0 & \text{, if } x < F_Z^{-1}(\tau) \end{cases}$$

(Aka "quantile regression" loss)

Roger Koenker, Quantile Regression, 2005

13

# The Quantile Regression Loss

▸ Given that the derivative of $L(x; Z, \tau)$ is $F_Z(x) - \tau$, we can recover the QR loss by integration

Quantile regression (QR) loss:

$$L_{QR}(x; Z, \tau) = (\tau - 1)\int_{-\infty}^{x}(z - x)dF_Z(z) + \tau\int_{x}^{\infty}(z - x)dF_Z(z)$$

(It is easy to verify that $\dfrac{d}{dx}L_{QR}(x; Z, \tau) = F_Z(x) - \tau$ by the Leibniz integral rule)

Alternative expression of QR loss:

$$\frac{d}{dx}\left(\int_{a(x)}^{b(x)} f(x,t)\,dt\right)$$

$$= f(x, b(x)) \cdot \frac{d}{dx}b(x) - f(x, a(x)) \cdot \frac{d}{dx}a(x) + \int_{a(x)}^{b(x)} \frac{\partial}{\partial x}f(x,t)\,dt$$

$$\rho_\tau(y) := y(\tau - \mathbb{I}\{y < 0\})$$

$$L_{QR}(x; Z, \tau) = E_Z[\rho_\tau(Z - x)]$$

14

# Summary: Loss Function of QR-DQN

$$L_{QRDQN}(s, a, r, s'; \theta) := \sum_{i=1}^{N} L_{QR}(B^*Z_{\bar{\theta}}(s, a); Z_{\theta}(s, a), \tau_i)$$

$$= \sum_{i=1}^{N} \mathbb{E}_{z \sim B^*Z_{\bar{\theta}}(s,a)}[\rho_{\tau_i}(z - Z_{\theta}(s, a))]$$
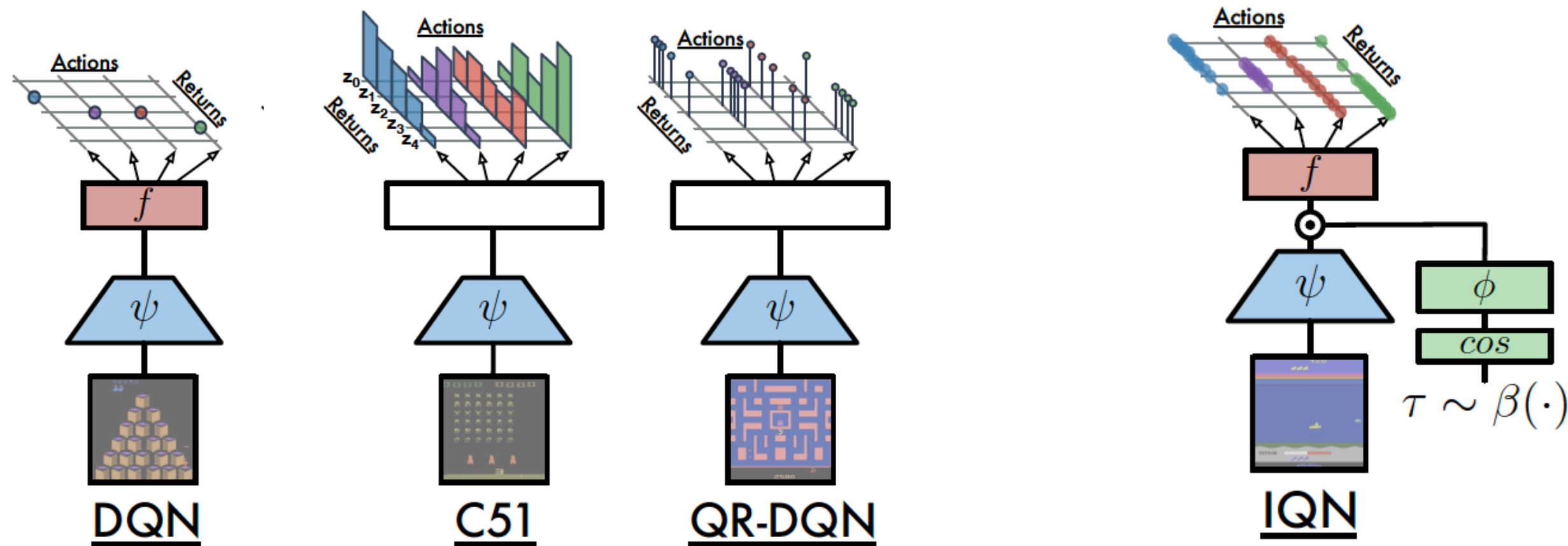
*N different quantiles*

▸ **Question**: Is $L_{QRDQN}(s, a, r, s'; \theta)$ easy to compute during training?

# Implicit Quantile Networks (IQN)

Dabney et al., Implicit Quantile Networks for Distributional Reinforcement Learning, ICML 2018

# IQN: A Generative Approach to Distributional RL

▸ An illustrative comparison of distributional Q-learning methods



DQN      C51      QR-DQN      IQN

Distributional RL via explicitly expressing the distribution $Z(s, a)$

Distributional RL via a generative model for distribution $Z(s, a)$

➡ Need sufficiently many atoms or quantiles for an accurate representation of $Z(s, a)$

(Figure credit: Will Dabney)

# Calculate QR Loss by *Sampling*
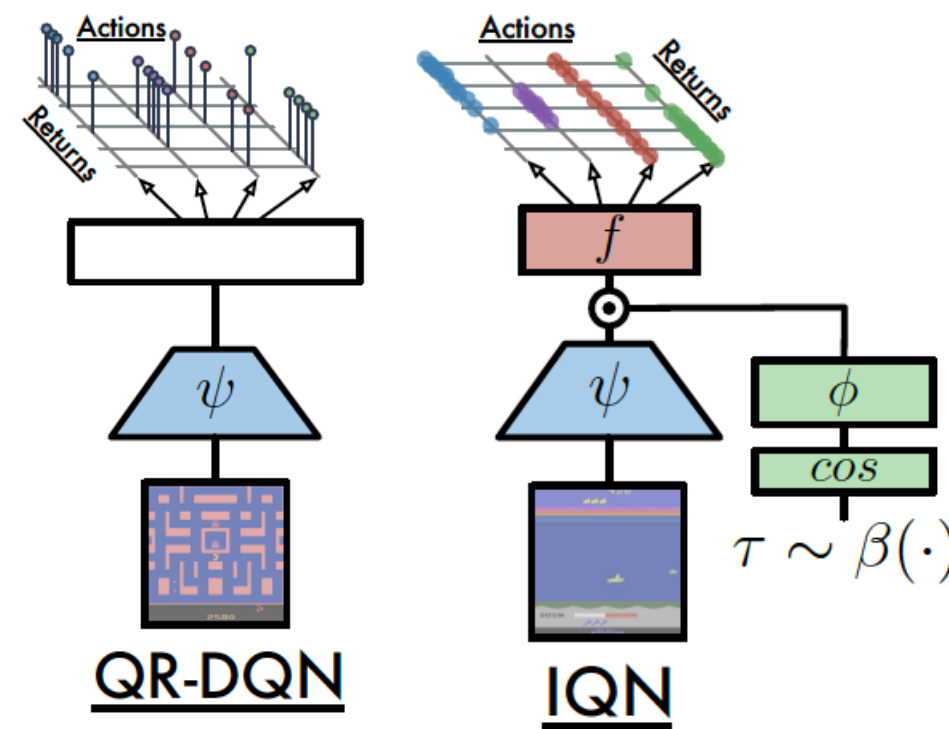
**QR loss:**

$$\rho_\tau(y) := y(\tau - \mathbb{1}\{y < 0\})$$
$$L_{QR}(x; Z, \tau) = E_{z \sim Z}[\rho_\tau(z - x)]$$



QR-DQN    IQN

▸ Recall QR-DQN:

   ▸ The QR loss is calculated explicitly

   ▸ $Z \Rightarrow$ target distribution induced by $\{\bar{\theta}_1, \cdots, \bar{\theta}_N\}$

▸ Idea: A practical way to calculate the QR loss is *sampling*!

$$L_{QR}(x; Z, \tau) \approx$$

▸ IQN **implicitly** parameterizes $Z$ by constructing a generator for $Z$

Suppose we are given a distribution of $Z$, denoted by $F_Z$ (CDF).

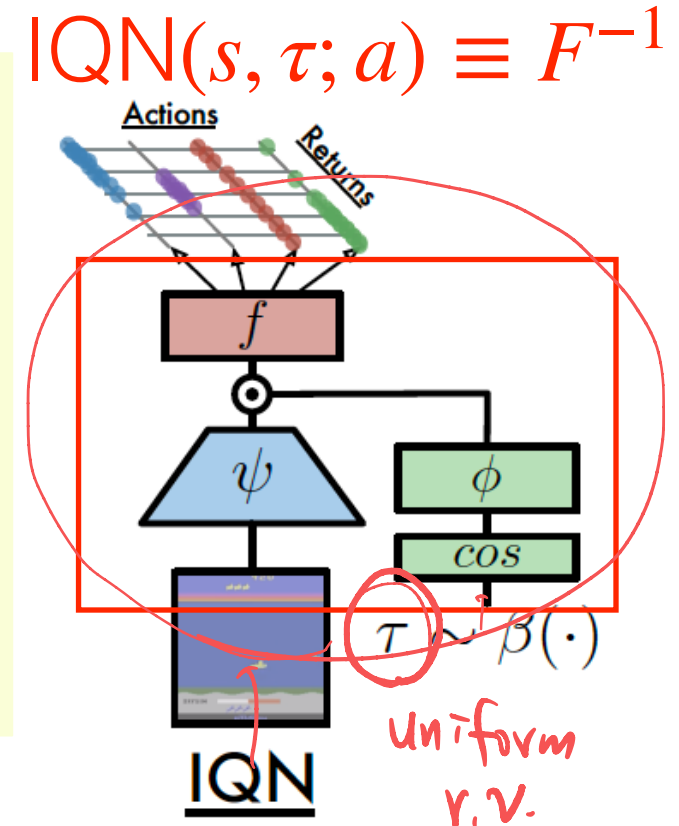Q: How to generate "random variables" of CDF $F_Z$ ?

I T S
↓ ↓ ↘ sampling
inverse transform

# QR Loss and Inverse Transform Sampling

**QR loss:**

$$\rho_\tau(y) := y(\tau - \mathbb{I}\{y < 0\})$$

$$L(x; Z, \tau) = E_{z \sim Z}[\rho_\tau(z - x)] \approx \frac{1}{K} \sum_{k=1}^{K} \rho_\tau(z_k - x)$$

$$(z_1, \cdots, z_K \sim Z)$$

$$\text{IQN}(s, \tau; a) \equiv F^{-1}$$

$$\tau \sim \beta(\cdot)$$

Uniform r.v.

**Inverse Transform Sampling (ITS)**: Generate <u>any</u> random variable with CDF $F$ from a uniform random variable

1. Generate a random variable $U \sim \text{Unif}(0,1)$

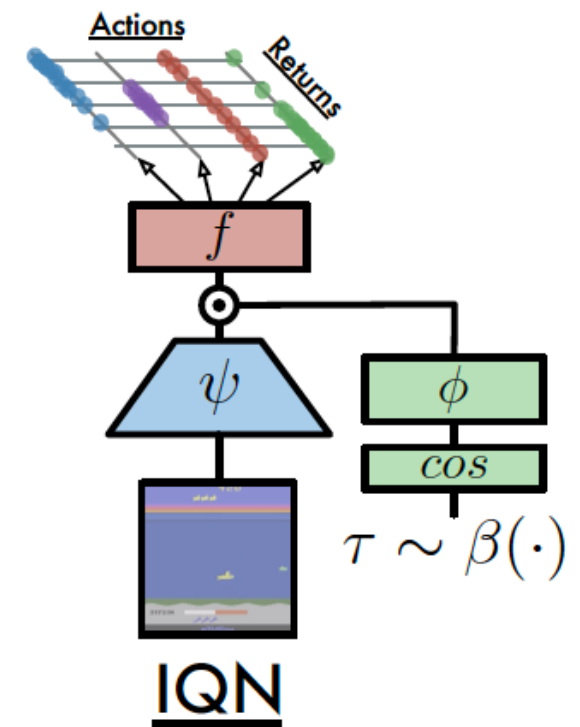2. Let $X = F^{-1}(U)$, where $F^{-1}(u) := \inf\{z : F(z) \geq u\}$

‣ ITS is essentially a generative approach!

# Calculating QR Loss in IQN

**QR loss:**

$$\rho_\tau(y) := y(\tau - \mathbb{I}\{y < 0\})$$

$$L(x; Z, \tau) = E_{z \sim Z}[\rho_\tau(z - x)] \approx \frac{1}{K} \sum_{k=1}^{K} \rho_\tau(z_k - x)$$

$$(z_1, \cdots, z_K \sim Z)$$



IQN

(Recall that $Z$ corresponds to the target distribution in QR-DQN)

At each update, given $(s, a, r, s')$, for a given $\tau \in [0,1]$:

  1. Draw $\tau'_1, \cdots, \tau'_K \sim \text{Unif}(0,1)$ ← a generative step!

  2. Get $z_1, \cdots, z_K$ by $z_i = r + \gamma \cdot \overline{\text{IQN}}(s', a'; \tau'_i)$

  3. QR loss in IQN $= \frac{1}{K} \sum_{i=1}^{K} \rho_\tau(z_i - \text{IQN}(s, a; \tau))$

(can be readily extended to multiple $\tau$)

# IQN is closely related to the <u>reparameterization trick</u>

- Suppose we want to compute a loss $L(\theta) = E_{X \sim p_\theta}[f(X)]$

  - $X$ is a random variable, and $p_\theta$ is the underlying distribution of $X$

- Question: $\nabla_\theta L(\theta) = ?$

$$\nabla_\theta L(\theta) = \nabla_\theta E_{X \sim p_\theta}[f(X)] = \nabla_\theta \left( \int f(x) p_\theta(x) dx \right)$$

$$= \int \left( f(x) \frac{1}{p_\theta(x)} \nabla_\theta p_\theta(x) \right) p_\theta(x) dx$$

$$= \int \left( f(x) \underline{\nabla_\theta \log p_\theta(x)} \right) p_\theta(x) dx$$

Easy to evaluate?

- Reparameterization trick: $\varepsilon \sim p(\varepsilon), \ L(\theta) = E_{\varepsilon \sim p}[g_\theta(\varepsilon)]$

$$\nabla_\theta L(\theta) = \nabla_\theta E_{\varepsilon \sim p}[g_\theta(\varepsilon)] = E_{\varepsilon \sim p}[\nabla_\theta g_\theta(\varepsilon)] \approx \frac{1}{K} \sum_{i=1}^{K} \nabla_\theta g_\theta(\varepsilon_i)$$

$$(\varepsilon_1, \cdots, \varepsilon_K \sim p)$$

Kingma and Welling, Auto-Encoding Variational Bayes, ICLR 2014