# **Distributional Bellman Operators over Mean Embeddings**

Li Kevin Wenliang, Grégoire Delétang, Matthew Aitchison, Marcus Hutter, Anian Ruoss, Arthur Gretton, Mark Rowland

- A new <u>family of distributional Bellman operators</u> for distributional reinforcement learning
- Return distributions encoded as <u>mean embeddings</u>, Bellman updates performed in ME space
- The algorithm is simple, fast and elegant
- Convergence by <u>error analysis</u>, nice Atari suite result

## **Distributional Bellman Consistency**

Consider policy evaluation over a Markov reward process

 $x \in \mathcal{X}, a \in \mathcal{A}, P_X(X|x,a), P_R(R|x,a), \gamma$ discount reward transition states actions

Trajectory under policy  $\pi$ :  $X_0 = x$ ,  $A_0$ ,  $R_0$ ,  $X_1$ ,  $A_1$ ,  $R_1$ , ...

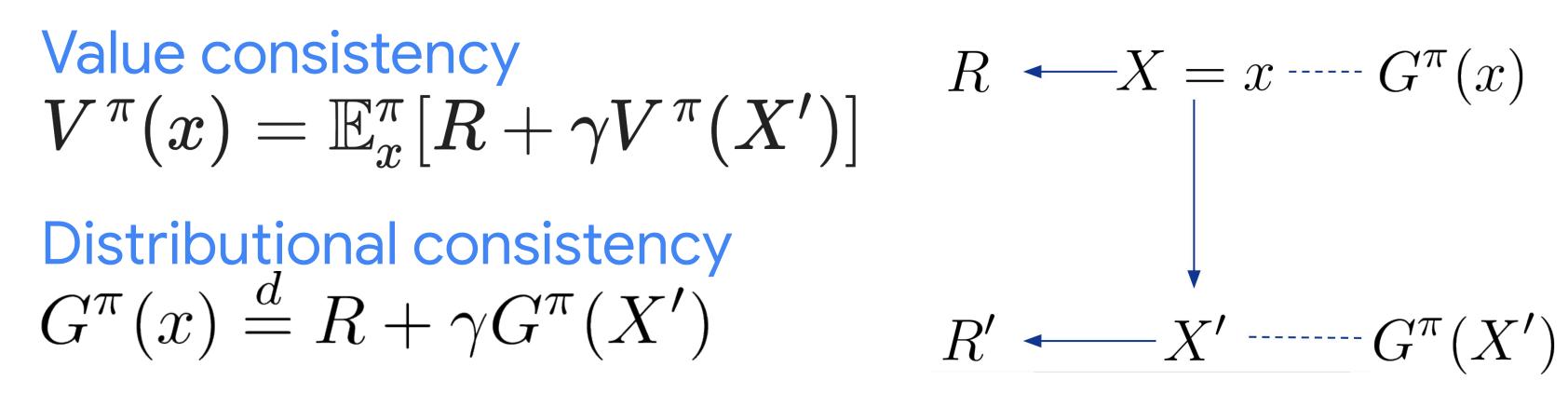
In value-based RL, we estimate the value for each state

$$V^{\pi}(x) := \mathbb{E}_x^{\pi}[G^{\pi}(x)], \quad G^{\pi}(x) := \sum_{t=0}^{\infty} \gamma^t R_t \mid X_0 = x$$

In distributional RL we estimate the return distribution

$$\eta^{\pi}(x) := \operatorname{Law}(G^{\pi}(x)) \quad \Leftrightarrow \quad G^{\pi}(x) \sim \eta^{\pi}(x)$$

Bellman consistencies: for a transition x, A, R, X'

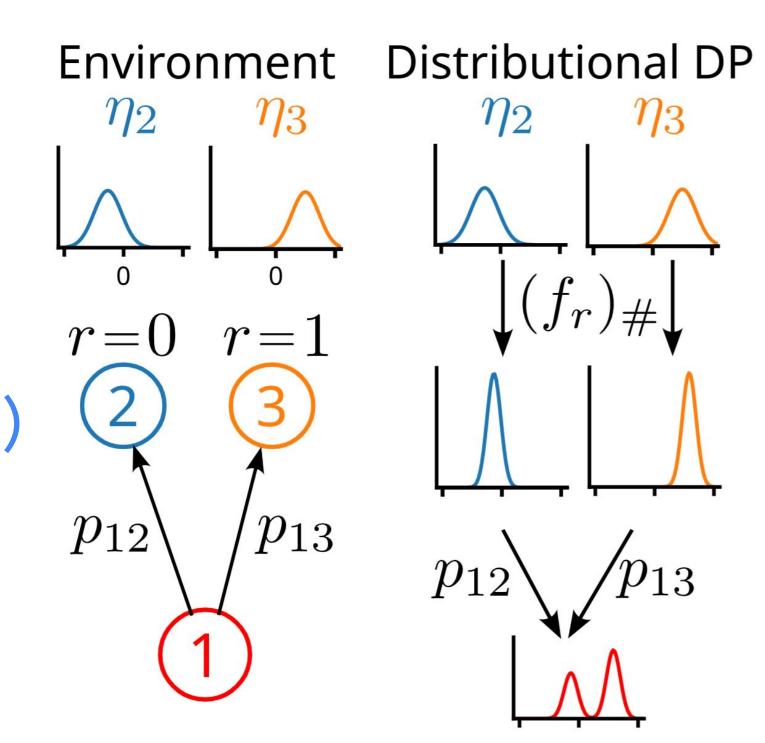


## **Distributional Bellman Update**

Dynamical programming (DP)  $V(x) \leftarrow \mathbb{E}^\pi_x[R + \gamma V(X')]$ Easy to implement with vectors Distributional DP (Bellemare+ 23)

$$G(x) \stackrel{?}{\leftarrow} R + \gamma G(X')$$

How to implement?



#### Main Idea: Bellman + Mean Embeddings

The mean embedding (ME) encodes distributions by feature functions (statistical functional or sketch):

 $ext{ for } egin{cases} \eta \in \mathscr{P}(\mathbb{R}), \ \phi: \mathbb{R} o \mathbb{R}^m, \ ext{ define } u := \mathbb{E}_{Z \sim \eta}[\phi(Z)]\,. \end{cases}$ 

Let's write the Bellman consistency using ME

$$U^\pi(x):=\mathbb{E}^\pi_x[\phi(G^\pi(x))]=\mathbb{E}^\pi_x[\phi(R+\gamma G^\pi(X'))]\,.$$

If the feature function satisfies:

$$\phi(r+\gamma z)=B_r\phi(z)$$
 for some  $B_r\in \mathbb{R}^{m imes m}$ , (2)

then we have a "closed" expression:

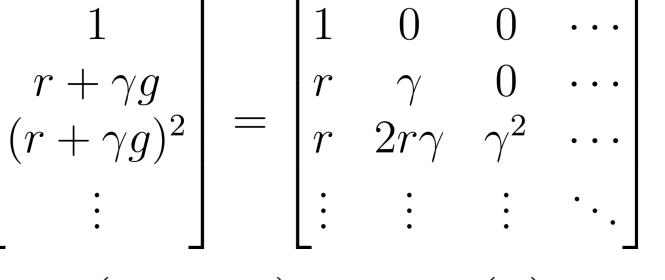
 $U^\pi(x) = \mathbb{E}^\pi_x[\phi(R+\gamma G^\pi(X'))] = \mathbb{E}^\pi_x[B_R\phi(G^\pi(X'))]$  $=\mathbb{E}_x^\pi[B_R U^\pi(X')]$  .

Leading to Sketch-DP  $U(x) \leftarrow \mathbb{E}_x^{\pi}[B_R U(X')]$ 

### On the Big If...

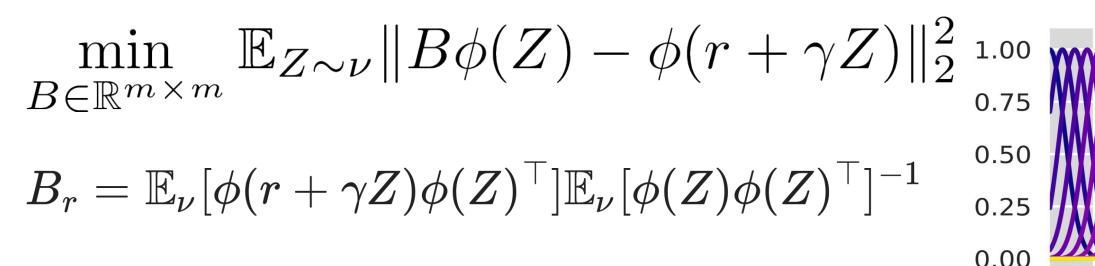
Does the assumption in (2) hold?

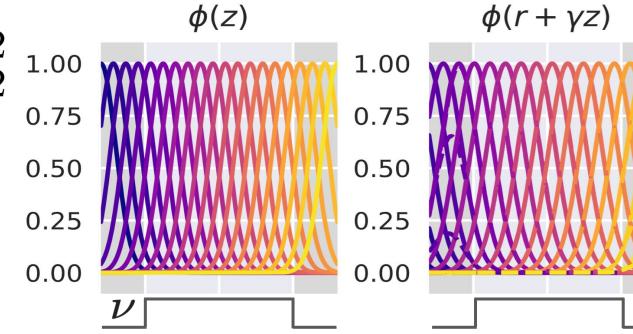
Rowland+ (2019) showed that only polynomials satisfy this assumption.

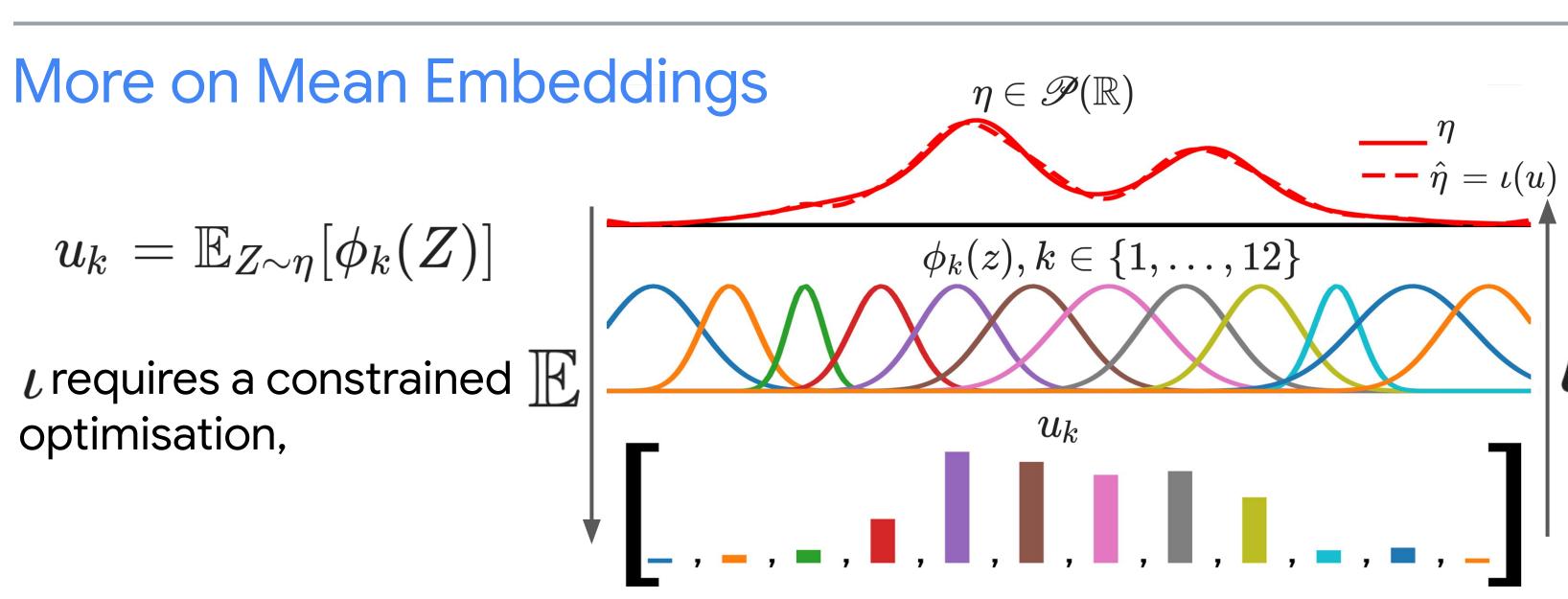


BUT let's relax (2) to an approximation:  $\phi(r + \gamma z) \approx B_r \phi(z)$ 

To ensure small error, solve linear regression for each r

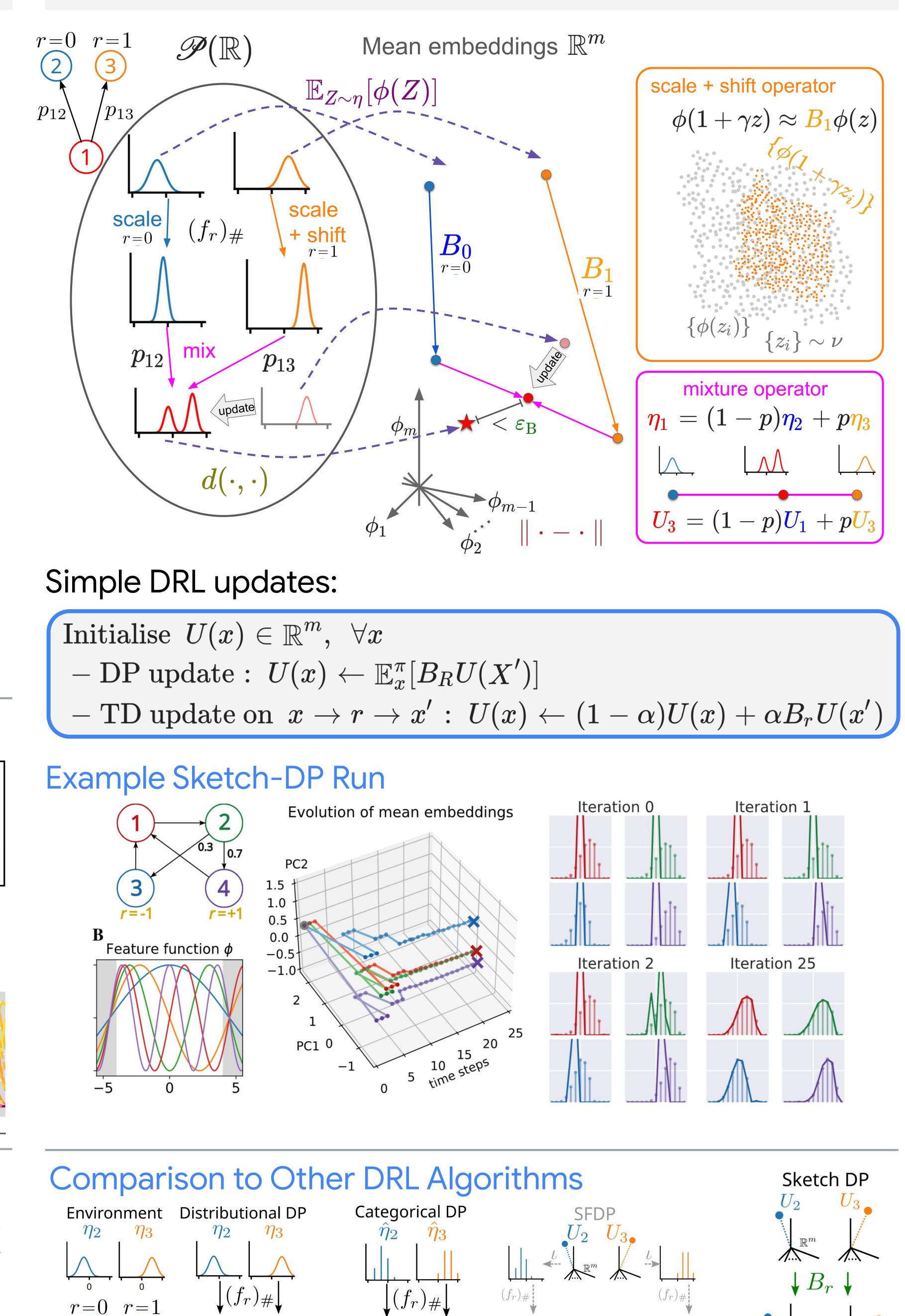








**Algorithms: Sketch-DP and Sketch-TD** 



┕╾╾╾┶

 $p_{12}$ 

 $p_{12}$ 

<del>┟╴╴┙╹┇╹╶╸┍╶╹╵╹╹╹</del>

 $p_{12} | / p_{13}$ 

╷┈╢╢╷──┝│

Google DeepMind





### **Convergence: Error Analysis**

Intuitively, to ensure convergence to  $U^{\pi}$ , we want  $\phi$  to be "rich":

- . approximate shift+scale versions well;
- 2. retain sufficient distributional information.

**Informal theorem:** Sketch-DP iterates  $\rightarrow$  neighbourhood of  $U^{\pi}$  if:

1: Feature approx. error bounds distributional update error

 $\|\mathbb{E}^\pi_x[\phi(R+\gamma G(X'))]-\mathbb{E}^\pi_x[B_RU(X')]\|$  $g \in [G_{\min}, G_{\max}]$  $=: \varepsilon_{\mathrm{R}}$ 

2: U and U' of two return distributions  $\eta$  and  $\eta'$  must satisfy:

- . there exists a distributional metric d on which the distributional Bellman operator contracts;
- $\mathsf{b.} \ \exists \, \varepsilon_1, \varepsilon_2 > 0 \ \text{ s.t. } \ \boldsymbol{d}(\eta, \eta') \varepsilon_1 \leq \|U U'\| \leq \boldsymbol{d}(\eta, \eta') + \varepsilon_2 \, .$

See paper for a concrete example using a particular sketch.

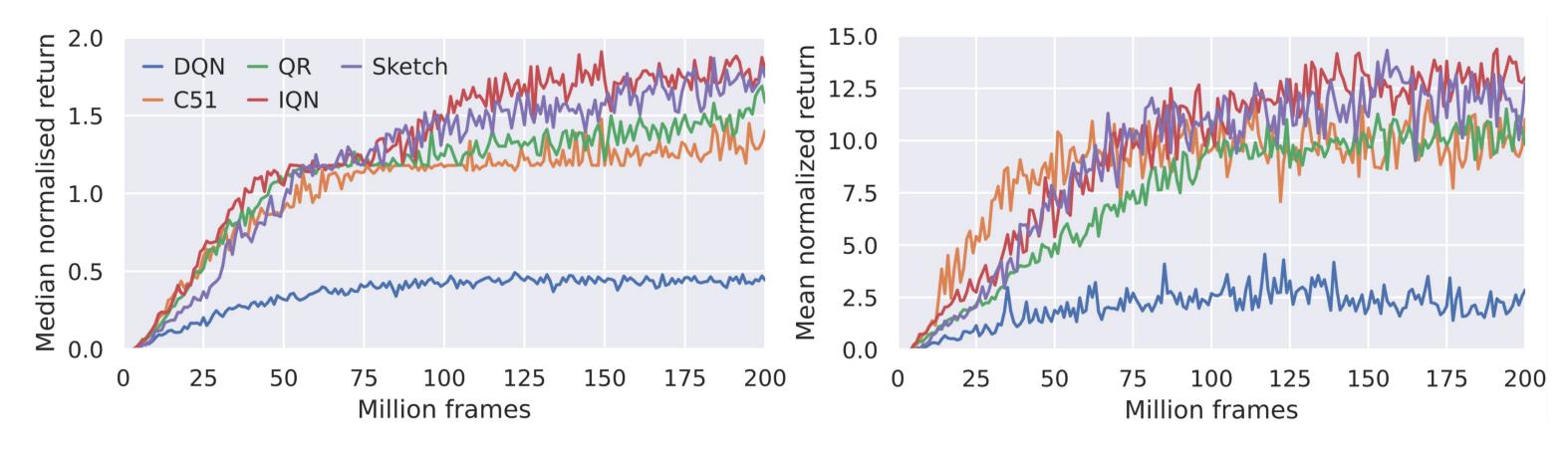
# Experiments

#### Main Atari Suite

ɛ-greedy policy, value readout  $\,Q(x,a)=eta\cdot U(x,a)$ 

 $\arg\min_{\beta} \mathbb{E}_{G \sim \mu} [(G - \langle \beta, \phi(G) \rangle)^2]$ 

sigmoid features  $\phi$ , based on QR-DQN architecture



#### Synthetic Experiments

 $p_{12} / p_{13}$ 

A few simple environments (column), feature types (hue) Cramer distance  $\max \ell_2[\iota(U(x)), \eta^{\pi}(x)]$ 

