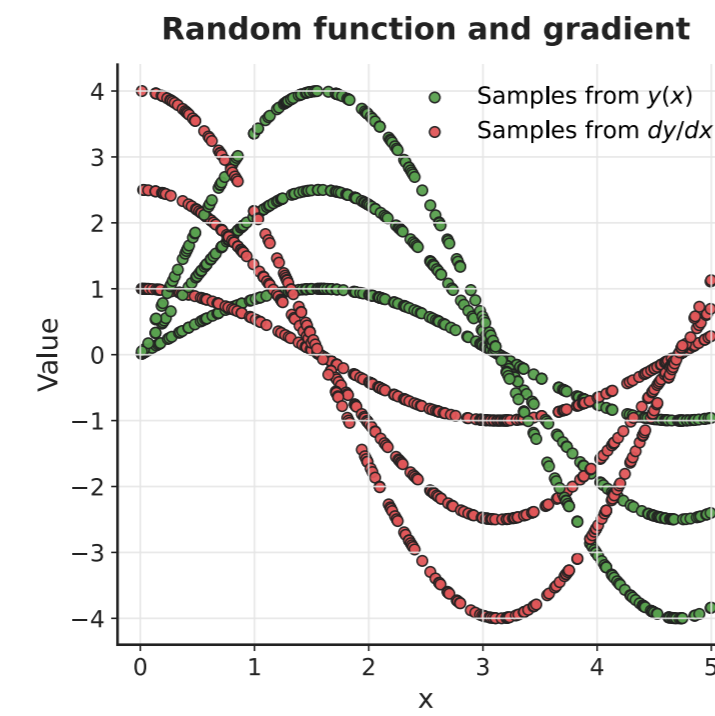




TL;DR

- Observations
 - Stochastic returns
 - Value gradients are stochastic
 - Deterministic value gradients may struggle
- Goals
 - Robustness & sample efficiency
- Contributions
 - Distributional Value Gradients
 - Joint law over return and gradient
 - Sobolev TD framework
 - Contraction theory
 - Trade-off interpretation



Value Gradients

- Policy optimization
 - Actor update from $\nabla_a Q^\pi(s, a)$

$$\nabla_\theta J(\theta) \propto \nabla_\theta \pi_\theta(s) \nabla_a Q^\pi(s, a)|_{a=\pi_\theta(s)}$$
- Learned dynamics
 - Gradient through one-step model

$$\hat{s}' = f(s, a), \quad \hat{r} = r(s, a)$$

$$\hat{\delta}(s, a) = \hat{r} + \gamma Q_{\phi'}(\hat{s}', \pi(\hat{s}'))$$
- TD on value and gradient

$$\mathcal{L}(\phi) = \|Q_\phi(s, a) - \hat{\delta}(s, a)\|^2 + \lambda \|\nabla_a Q_\phi(s, a) - \nabla_a \hat{\delta}(s, a)\|^2$$

Distributional RL

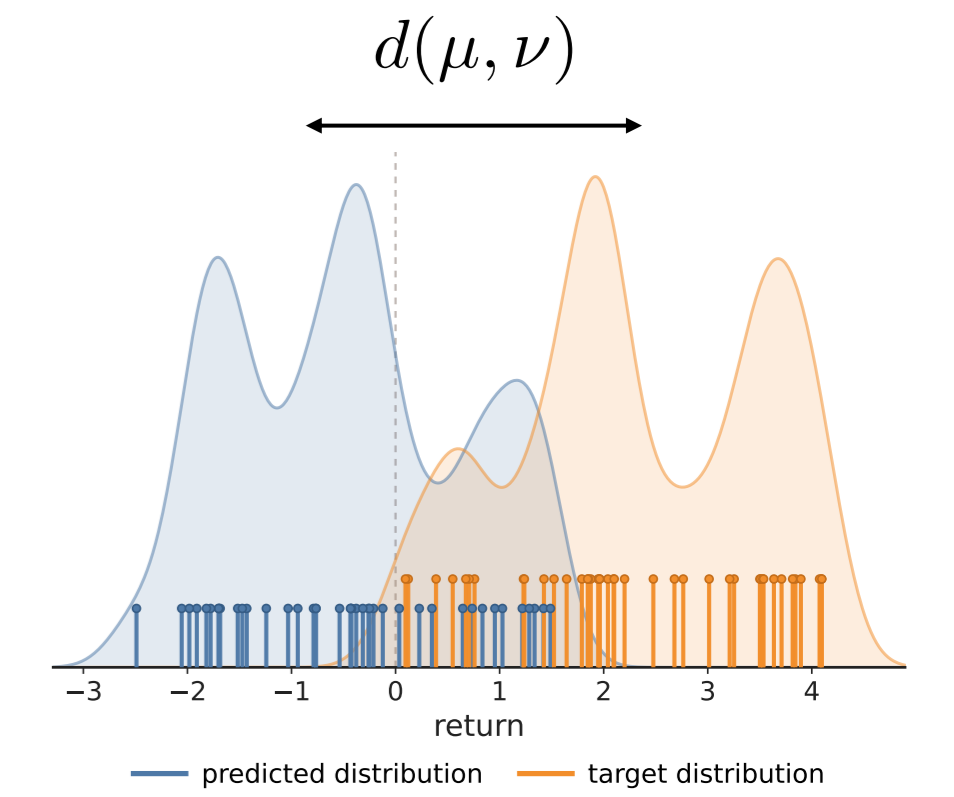
- Distribution over returns

$$\eta^\pi(s, a) = \text{Law}[Z^\pi(s, a)]$$

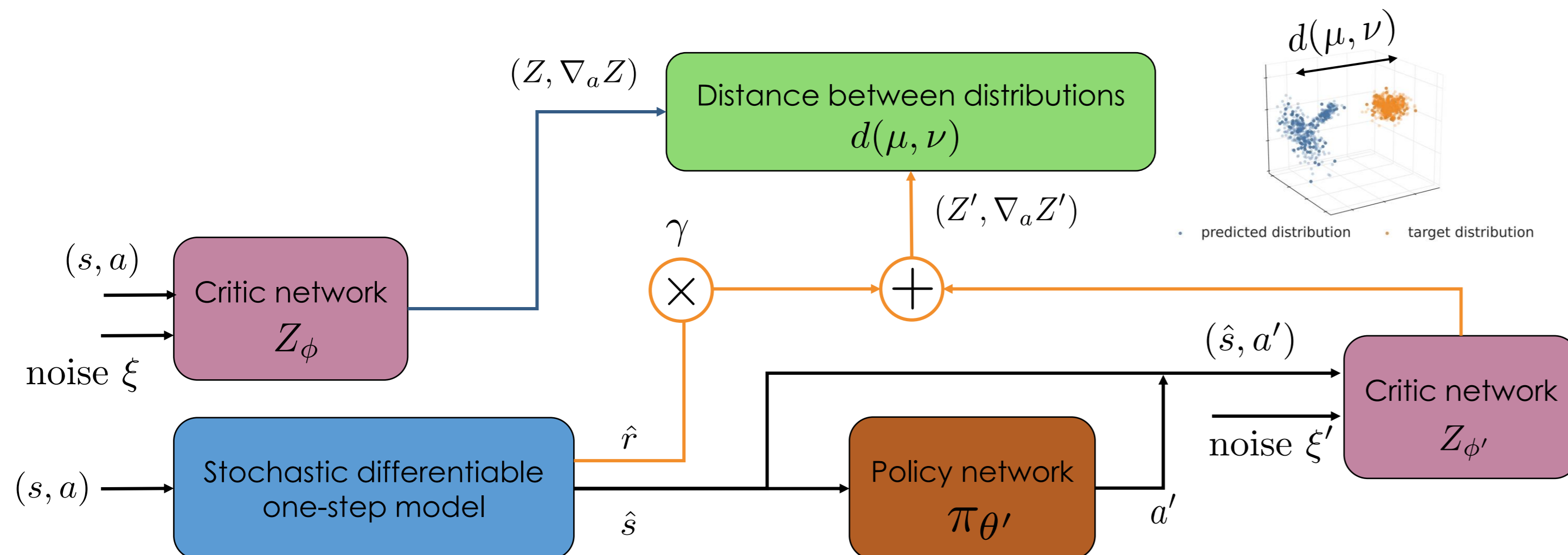
$$\eta_{\text{tgt}}(s, a) = \text{Law}[r + \gamma Z_{\phi'}(s', \pi(s'))]$$
- Distance between distributions

$$d(\mu, \nu)$$
 i.e. W_p , MMD
- TD on distributions

$$\mathcal{L}(\phi) = d(\eta_\phi(s, a), \eta_{\text{tgt}}(s, a))$$



Distributional Value Gradients

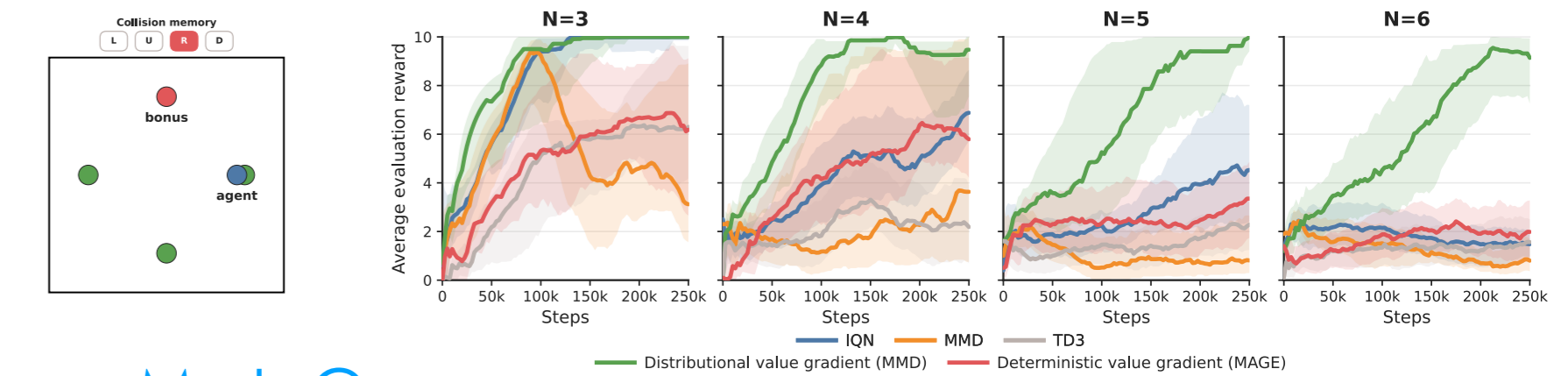


- TD on joint distribution

$$(Z_\phi(s, a, \xi), \nabla_a Z_\phi(s, a, \xi)), \quad \xi \sim \mathcal{N}(0, I)$$
- Stochastic world model
 - Environment aleatoric uncertainty
 - Cheap sampling and gradients
- Theory
 - Contraction in Wasserstein $\gamma \kappa < 1$
 - Smoothness κ
 - Trade-off
 - Large physical gradients
 - Lower horizon γ

Results

Toy environment



MuJoCo

