

409 A KL-Divergence Trust Region Projection Layer

410 As already mentioned in the main text, TRPLs [9] present a scalable and mathematically sound
 411 approach for enforcing trust regions in step-based deep RL. The layer takes the output of a standard
 412 Gaussian policy as input in terms of mean μ and variance Σ and projects it into the trust region if
 413 the given mean and variance violate their respective bounds. This projection is done for each input
 414 state individually. Subsequently, the projected Gaussian policy distribution with parameters $\tilde{\mu}$, $\tilde{\Sigma}$ is
 415 used for any further steps, e. g. for sampling and/or loss computation. Formally, the layer solves the
 416 following two optimization problems for each state s

$$\arg \min_{\tilde{\mu}_s} d_{\text{mean}}(\tilde{\mu}_s, \mu(s)), \quad \text{s.t.} \quad d_{\text{mean}}(\tilde{\mu}_s, \mu_{\text{old}}(s)) \leq \epsilon_{\mu}, \quad \text{and} \quad (1)$$

$$\arg \min_{\tilde{\Sigma}_s} d_{\text{cov}}(\tilde{\Sigma}_s, \Sigma(s)), \quad \text{s.t.} \quad d_{\text{cov}}(\tilde{\Sigma}_s, \Sigma_{\text{old}}(s)) \leq \epsilon_{\Sigma}, \quad (2)$$

417 where $\tilde{\mu}_s$ and $\tilde{\Sigma}_s$ are the optimization variables for input state s and ϵ_{μ} and ϵ_{Σ} are the trust region
 418 bounds for mean and covariance, respectively. Finally, μ_{old} and Σ_{old} are the reference mean and
 419 covariance for the trust region and d_{mean} as well as d_{cov} are the similarity metrics for the mean
 420 and covariance of a decomposable distance or divergence measure. As we only leverage the KL-
 421 divergence projection, we will provide only details for this particular projection below. For the other
 422 two projections we refer the reader to Otto et al. [9].

423 Inserting the mean part of the Gaussian KL divergence into Equation 1 yields

$$\arg \min_{\tilde{\mu}} (\mu - \tilde{\mu})^T \Sigma_{\text{old}}^{-1} (\mu - \tilde{\mu}) \quad \text{s.t.} \quad (\mu_{\text{old}} - \tilde{\mu})^T \Sigma_{\text{old}}^{-1} (\mu_{\text{old}} - \tilde{\mu}) \leq \epsilon_{\mu}.$$

424 After differentiating the dual w.r.t. $\tilde{\mu}$, we can solve for the projected mean

$$\tilde{\mu} = \frac{\mu + \omega \mu_{\text{old}}}{1 + \omega} \quad \text{with} \quad \omega = \sqrt{\frac{(\mu_{\text{old}} - \mu)^T \Sigma_{\text{old}}^{-1} (\mu_{\text{old}} - \mu)}{\epsilon_{\mu}}} - 1,$$

425 leveraging the optimal Lagrange multiplier ω . Similarly, we can insert the covariance part of the
 426 Gaussian KL divergence into Equation 2, which results in

$$\arg \min_{\tilde{\Sigma}} \text{tr}(\Sigma^{-1} \tilde{\Sigma}) + \log \frac{|\Sigma|}{|\tilde{\Sigma}|}, \quad \text{s.t.} \quad \text{tr}(\Sigma_{\text{old}}^{-1} \tilde{\Sigma}) - d + \log \frac{|\Sigma_{\text{old}}|}{|\tilde{\Sigma}|} \leq \epsilon_{\Sigma},$$

427 where d is the number of degrees of freedom (DoF). Once again, differentiating and solving the dual
 428 $g(\eta)$ for the projected covariance yields

$$\tilde{\Sigma} = \left(\frac{\eta^* \Sigma_{\text{old}}^{-1} + \Sigma^{-1}}{\eta^* + 1} \right)^{-1} \quad \text{with} \quad \eta^* = \arg \min_{\eta} g(\eta), \quad \text{s.t.} \quad \eta \geq 0.$$

429 Here, the the optimal Lagrange multiplier η^* cannot be computed in closed form, however, a stan-
 430 dard numerical optimizer, such as BFGS, is able to efficiently find it. This can be made differentiable
 431 by taking the differentials of the KKT conditions of the dual. For more details, we refer to the orig-
 432 inal work [9].

433 B Environment Details

434 B.1 Box Pushing

435 The goal of the box pushing task is to move a box to a specified goal location and orientation using
 436 the seven DoF Franka Emika Panda. Hence, the context space for this task is the goal position
 437 $x \in [0.3, 0.6]$, $y \in [-0.45, 0.45]$ and the goal orientation $\theta \in [0, 2\pi]$. In addition to the contexts, the
 438 observation space for the step-based algorithms contains information about joints and end-effector
 439 as well as the current box location and orientation. To the original torque from the policy we add

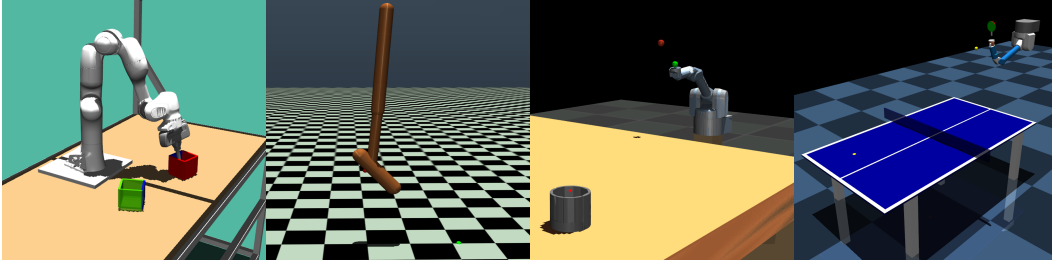


Figure 5: Visualization of the four control tasks box pushing, hopper jumping, beer pong, and table tennis.

440 gravity compensation in each time step. The task is considered successfully solved if the position
 441 distance $\leq 0.05\text{m}$ and the orientation error $\leq 0.5\text{rad}$. For the total reward we consider different
 442 sub-rewards. First, the distance to the goal

$$R_{\text{goal}} = \|\mathbf{p} - \mathbf{p}_{\text{goal}}\|,$$

443 where \mathbf{p} is the box position and \mathbf{p}_{goal} the goal position itself. Second, the rotation distance

$$R_{\text{rotation}} = \frac{1}{\pi} \arccos |\mathbf{r} \cdot \mathbf{r}_{\text{goal}}|,$$

444 where \mathbf{r} and \mathbf{r}_{goal} are the box orientation and goal orientation in quaternion, respectively. Third, an
 445 incentive to keep the rod within the box

$$R_{\text{rod}} = \text{clip}(\|\mathbf{p} - \mathbf{h}_{\text{pos}}\|, 0.05, 10)$$

446 where \mathbf{h}_{pos} is the position of the rod tip. Fourth, a similar incentive that encourages to maintain the
 447 rod in a desired rotation

$$R_{\text{rod.rotation}} = \text{clip}\left(\frac{2}{\pi} \arccos |\mathbf{h}_{\text{rot}} \cdot \mathbf{h}_0|, 0.25, 2\right),$$

448 where \mathbf{h}_{rot} and $\mathbf{h}_0 = (0.0, 1.0, 0.0, 0.0)$ are the current and desired rod orientation in quaternion,
 449 respectively. And lastly, we utilize the following error

$$\text{err}(\mathbf{q}, \dot{\mathbf{q}}) = \sum_{i \in \{i \mid |q_i| > |q_i^b|\}} (|q_i| - |q_i^b|) + \sum_{j \in \{j \mid |\dot{q}_j| > |\dot{q}_j^b|\}} (|\dot{q}_j| - |\dot{q}_j^b|).$$

450 Here, \mathbf{q} , $\dot{\mathbf{q}}$, \mathbf{q}^b , and $\dot{\mathbf{q}}^b$ are the robot joint's position and velocity as well as their respective bounds.

451 Additionally, we consider an action cost in each time step t

$$\tau_t = 5 \cdot 10^{-4} \sum_i^K (a_t^i)^2,$$

452 where $K = 7$ is the number of DoF. In total we consider three different rewards.

453 **Dense Reward.** The dense reward provides information about the goal and rotation distance in each
 454 time step t on top of the utility rewards

$$R_{\text{tot}} = -R_{\text{rod}} - R_{\text{rod.rotation}} - \tau_t - \text{err}(\mathbf{q}, \dot{\mathbf{q}}) - 3.5R_{\text{goal}} - 2R_{\text{rotation}}.$$

455 **Time-Dependent Sparse Reward.** The time-dependent sparse reward is similar to the dense re-
 456 ward, but only returns the goal and rotation distance in the last time step T

$$R_{\text{tot}} = \begin{cases} -R_{\text{rod}} - R_{\text{rod.rotation}} - \tau_t - \text{err}(\mathbf{q}, \dot{\mathbf{q}}), & t < T, \\ -R_{\text{rod}} - R_{\text{rod.rotation}} - \tau_t - \text{err}(\mathbf{q}, \dot{\mathbf{q}}) - 350R_{\text{goal}} - 200R_{\text{rotation}}, & t = T. \end{cases}$$

457 **Time- and Space-Dependent Sparse Reward.** The second sparse reward additionally adds sparsity
 458 based on the position and only returns goal and rotation distance in the last time step when the box
 459 is near the goal location

$$R_{\text{tot}} = \begin{cases} -R_{\text{rod}} - R_{\text{rod.rotation}} - \tau_t - \text{err}(\mathbf{q}, \dot{\mathbf{q}}) \cdots \\ \cdots - \text{clip}(1050R_{\text{goal}}, 0, 100) - \text{clip}(15R_{\text{rotation}}, 0, 100) + 300, & t = T \text{ and } R_{\text{goal}} \leq 0.1, \\ -R_{\text{rod}} - R_{\text{rod.rotation}} - \tau_t - \text{err}(\mathbf{q}, \dot{\mathbf{q}}), & \text{else.} \end{cases}$$

460 **B.2 Hopper Jump**

461 In the hopper jump task the agent has to learn to jump as high as possible and land on a certain goal
 462 position at the same time. We consider five basis functions per joint resulting in an 15 dimensional
 463 weight space. The context is four-dimensional consisting of the initial joint angles $\theta \in [-0.5, 0]$, $\gamma \in$
 464 $[-0.2, 0]$, $\phi \in [0, 0.785]$ and the goal landing position $x \in [0.3, 1.35]$. We consider a non-Markovian
 465 reward function for the episode-based algorithms and a step-based reward for PPO, which we have
 466 extensively designed to obtain the highest possible jump.

467 **Non-Markovian Reward.** In each time-step t we provide an action cost

$$\tau_t = 10^{-3} \sum_i^K (a_t^i)^2,$$

468 where $K = 3$ is the number of DoF. In the last time-step T of the episode we provide a reward
 469 which contains information about the whole episode as

$$\begin{aligned} R_{height} &= 10h_{max}, \\ R_{gdist} &= \|p_{foot,T} - p_{goal}\|_2, \\ R_{cdist} &= \|p_{foot,contact} - p_{goal}\|_2, \\ R_{healthy} &= \begin{cases} 2 & \text{if } z_T \in [0.5, \infty] \text{ and } \theta, \gamma, \phi \in [-\infty, \infty] \\ 0 & \text{else} \end{cases}, \end{aligned}$$

470 where h_{max} is the maximum jump height in z-direction of the center of mass reached during the
 471 whole episode, $p_{foot,t}$ is the x-y-z position of the foot’s heel at time step t , $p_{foot,contact}$ is the foot’s
 472 heel position when having a contact with the ground after the first jump, p_{goal} is the goal landing
 473 position of the heel. $R_{healthy}$ is a slightly modified reward of the healthy reward defined in the
 474 original hopper task. The hopper is considered as ‘healthy’ if the z position of the center of mass is
 475 within the range $[0.5m, \infty]$. This encourages the hopper to stand at the end of the episode. Note that
 476 all states need to be within the range $[-100, 100]$ for $R_{healthy}$. Since this is defined in the hopper
 477 task from OpenAI already, we haven’t mentioned it here. The total reward at the end of an episode
 478 is given as

$$R_{tot} = - \sum_{t=0}^T \tau_t + R_{height} + R_{gdist} + R_{cdist} + R_{healthy}.$$

479 **Step-Based Reward.** We consider a step-based alternative reward such that PPO is also able to
 480 learn a meaningful behavior on this task. We have tuned the reward such that we can obtain the
 481 best performance. The observation space is the same as in the original hopper task from OpenAI
 482 extended with the goal landing position and the current distance of the foot’s heel and the goal
 483 landing position. We again consider the action cost in each time-step t

$$\tau_t = 10^{-3} \sum_i^K (a_t^i)^2,$$

484 and additionally consider the rewards

$$\begin{aligned} R_{height,t} &= 3h_t \\ R_{gdist,t} &= 3\|p_{foot,t} - p_{goal}\|_2 R_{healthy,t} = \begin{cases} 1 & \text{if } z_t \in [0.5, \infty] \text{ and } \theta, \gamma, \phi \in [-\infty, \infty] \\ 0 & \text{else} \end{cases}, \end{aligned}$$

485 where these rewards are now returned to the agent in each time-step t , resulting in the reward per
 486 time-step

$$r_t(s_t, a_t) = -\tau_t + R_{height,t} + R_{gdist,t} + R_{healthy,t}.$$

487 **B.3 Beer Pong**

488 In the Beer Pong task the $K = 7$ Degrees of Freedom (DoF) robot has to throw a ball into a cup on
 489 a big table. The context is defined by the cup’s two dimensional position on the table which lies in
 490 the range $x \in [-1.42, 1.42]$, $y \in [-4.05, -1.25]$. For the step-based algorithms we consider cosine
 491 and sine of the robot’s angles, the angle velocities, the ball’s distance to the cup bottom, the ball’s
 492 distance to the cup’s top, the cup position and the current time step. The action space for the step-
 493 based algorithms is defined as the torques for each joint, the parameter space for the episode-based
 494 methods is 15 dimensional which consists of the two weights for the basis functions per joint and
 495 the duration of the throwing trajectory, i.e. the ball release time.

496 We generally consider action penalties

$$\tau_t = \frac{1}{K} \sum_i^K (a_t^i)^2,$$

497 consisting of the sum of squared torques per joint. For $t < T$ we consider the reward

$$r_t(s_t, a_t) = -\alpha_t \tau_t,$$

498 with $\alpha_t = 10^{-2}$. For $t = T$ we consider the non-Markovian reward

$$R_{task} = \begin{cases} -4 - \min(\|p_{c,top} - p_{b,1:T}\|_2^2) - 0.5\|p_{c,bottom} - p_{b,T}\|_2^2 \cdots & \text{if cond. 1} \\ \cdots - 2\|p_{c,bottom} - p_{b,k}\|_2^2 - \alpha_T \tau, & \text{if cond. 2} \\ -4 - \min(\|p_{c,top} - p_{b,1:T}\|_2^2) - 0.5\|p_{c,bottom} - p_{b,T}\|_2^2 - \alpha_T \tau, & \text{if cond. 3} \\ -2 - \min(\|p_{c,top} - p_{b,1:T}\|_2^2) - 0.5\|p_{c,bottom} - p_{b,T}\|_2^2 - \alpha_T \tau, & \text{if cond. 4} \\ -\|p_{c,bottom} - p_{b,T}\|_2^2 - \alpha_T \tau, & \end{cases}$$

$$R_{task} = \begin{cases} -4 - \min(\|p_{c,top} - p_{b,1:T}\|_2^2) - 0.5\|p_{c,bottom} - p_{b,T}\|_2^2 \cdots & \text{if cond. 1} \\ \cdots - 2\|p_{c,bottom} - p_{b,k}\|_2^2 - \alpha_T \tau, & \text{if cond. 2} \\ -4 - \min(\|p_{c,top} - p_{b,1:T}\|_2^2) - 0.5\|p_{c,bottom} - p_{b,T}\|_2^2 - \alpha_T \tau, & \text{if cond. 3} \\ -2 - \min(\|p_{c,top} - p_{b,1:T}\|_2^2) - 0.5\|p_{c,bottom} - p_{b,T}\|_2^2 - \alpha_T \tau, & \text{if cond. 4} \\ -\|p_{c,bottom} - p_{b,T}\|_2^2 - \alpha_T \tau, & \end{cases},$$

499 where $p_{c,top}$ is the position of the top edge of the cup, $p_{c,bottom}$ is the ground position of the cup,
 500 $p_{b,t}$ is the position of the ball at time point t , and τ is the squared mean torque over all joints during
 501 one rollout and $\alpha_T = 10^{-4}$. The different conditions are:

- 502 • cond. 1: The ball had a contact with the ground before having a contact with the table.
- 503 • cond. 2: The ball is not in the cup and had no table contact
- 504 • cond. 3: The ball is not in the cup and had table contact
- 505 • cond. 4: The ball is in the cup.

506 Note that $p_{b,k}$ is the ball’s and the ground’s contact position and is only given, if the ball had a
 507 contact with the ground first.

508 At time step $t = T$ we also give information whether the agent’s chosen ball release time B was
 509 reasonable

$$R_{release} = \begin{cases} -30 - 10(B - B_{min})^2, & \text{if } B < B_{min} \\ -30 - 10(B - B_{max})^2, & \text{if } B > B_{max} \end{cases},$$

510 where we define $B_{min} = 0.1s$ and $B_{max} = 1s$, such that the agent is encouraged to throw the ball
 511 within the time range $[B_{min}, B_{max}]$.

512 The total return over the whole episode is therefore given as

$$R_{tot} = \sum_{t=1}^{T-1} r_t(s_t, a_t) + R_{task} + R_{release}$$

513 A throw is considered as successful if the ball is in the cup at the end of an episode.

514 **B.4 Table Tennis**

515 We consider table tennis for the entire table, i. e. incoming balls are anywhere on the side of the robot
 516 and goal locations anywhere on the opponents side. The goal is to use the seven degree of freedom
 517 robotic arm to hit the incoming ball based on its landing position and return it as close as possible
 518 to the specified goal location. As context space we consider the initial ball position $x \in [-1, -0.2]$,
 519 $y \in [-0.65, 0.65]$ and the goal position $x \in [-1.2, -0.2]$, $y \in [-0.6, 0.6]$. The observation space
 520 again contains additional information about the joints and the ball. For this experiment, we do not
 521 use any gravity compensation and allow in the episode-based setting to learn the start time t_0 and the
 522 trajectory duration T . The task is considered successful if the returned ball lands on the opponent’s
 523 side of table and within $\leq 0.2\text{m}$ to the goal location. The reward is defined as

$$r_{task} = \begin{cases} 0, & \text{if cond. 1} \\ 0.2 - 0.2 \tanh(\min \|\mathbf{p}_r - \mathbf{p}_b\|^2), & \text{if cond. 2} \\ 3 - 2 \tanh(\min \|\mathbf{p}_r - \mathbf{p}_b\|^2) - \tanh(\|\mathbf{p}_l - \mathbf{p}_{goal}\|^2), & \text{if cond. 3} \\ 6 - 2 \tanh(\min \|\mathbf{p}_r - \mathbf{p}_b\|^2) - 4 \tanh(\|\mathbf{p}_l - \mathbf{p}_{goal}\|^2), & \text{if cond. 4} \\ 7 - 2 \tanh(\min \|\mathbf{p}_r - \mathbf{p}_b\|^2) - 4 \tanh(\|\mathbf{p}_l - \mathbf{p}_{goal}\|^2), & \text{if cond. 5} \end{cases}$$

524 where \mathbf{p}_r is the position of racket center, \mathbf{p}_b is the position of the ball, \mathbf{p}_l is the ball landing position,
 525 \mathbf{p}_{goal} is the target position. The different conditions are

- 526 • cond. 1: the end of episode is not reached
- 527 • cond. 2: robot did not hit the ball
- 528 • cond. 3: robot did hit the ball but the ball did not land on table or floor
- 529 • cond. 4: robot did hit the ball and returned it to the table or floor but it did not cross the net
- 530 • cond. 5: robot did hit the ball and returned it to the table or floor and cross the net

531 The episode ends when any of the following conditions are met

- 532 • the maximum horizon length is reached
- 533 • ball did land on the floor without hitting
- 534 • ball did land on the floor or table after hitting

535 For BBRL-PPO and BBRL-TRPL, the whole desired trajectory is obtained ahead of environment
 536 interaction, making use of this property we can collect some samples without physical simulation.
 537 The reward function based on this desired trajectory is defined as

$$r_{traj} = - \sum_{(i,j)} |\tau_{ij}^d| - |q_j^b|, \quad (i,j) \in \{(i,j) \mid |\tau_{ij}^d| > |q_j^b|\}$$

538 where τ^d is the desired trajectory, i is the time index, j is the joint index, q^b is the joint position
 539 upperbound. The desired trajectory is considered as invalid if $r_{traj} < 0$, an invalid trajectory will
 540 not be executed by robot. The overall reward for BBRL is defined as:

$$r = \begin{cases} r_{traj}, & r_{traj} < 0 \\ r_{task}, & \text{otherwise} \end{cases}$$

541 **C Additional Evaluations**

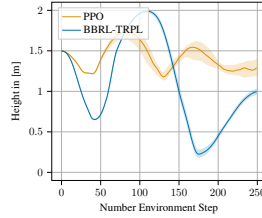


Figure 6: The improved performance on the Hopper Jump task is also demonstrated on the jumping profile for a fixed context. While BBRL-TRPL jumps once as high as possible, PPO constantly tries to maximize the height at each time step which leads to several jumps throughout the episode and consequently to a lower maximum height.

542 D Hyperparameters

543 For all methods we optimized the learning rate, sample size, batch size, number of layers, and the
 544 number of epochs. For all BBRL methods and NDP, we additionally optimized the number of basis
 545 functions. Moreover, we found that NDP requires tuning of the scale of the predicted DMP weights,
 546 which was hard-coded to 100 in the original code base. However, this value only worked for the
 547 meta-world tasks, but not for the other tasks, hence we adjusted it to allow for a fair comparison.

Table 1: Hyperparameters for the modified reacher experiments.

	PPO	NDP	BBRL-PPO	BBRL-TRPL
number samples	16000			64
GAE λ	0.95			n.a.
discount factor	0.99			n.a.
ϵ_μ		n.a.		0.05
ϵ_Σ		n.a.		0.0005
optimizer			adam	
epochs	10			100
learning rate			3e-4	
use critic	True			False
epochs critic	10			n.a.
learning rate critic	3e-4			n.a.
number minibatches	32			n.a.
trust region loss weight		n.a.		10.0
entropy loss penalty			0	
normalized observations	True			False
normalized rewards	True			False
observation clip	10.0			n.a.
reward clip	10.0			n.a.
critic clip		0.2		n.a.
importance ratio clip		0.2		n.a.
hidden layers			[32, 32]	
hidden layers critic	[32, 32]			n.a.
hidden activation			tanh	
initial std			1.0	
number basis functions	n.a.		5	
number zero basis	n.a.		1	
weight scale	n.a.	20		n.a.

Table 2: Hyperparameters for the box pushing experiments.

	PPO	NDP	BBRL-PPO	BBRL-TRPL
number samples	16000			160
GAE λ	0.95			n.a.
discount factor	0.99			n.a.
ϵ_μ			n.a.	0.005
ϵ_Σ			n.a.	0.0005
optimizer			adam	
epochs	10			100
learning rate	3e-4			1e-4
use critic	True			True
epochs critic	10			100
learning rate critic	3e-4			1e-4
number minibatches	40			n.a.
trust region loss weight		n.a.		25.0
entropy loss penalty			0	
normalized observations	True			False
normalized rewards	True			False
observation clip	10.0			n.a.
reward clip	10.0			n.a.
critic clip		0.2		n.a.
importance ratio clip		0.2		n.a.
hidden layers	[256, 256]			[128, 128]
hidden layers critic	[256, 256]			[32, 32]
hidden activation			tanh	
initial std	1.0			1.0
number basis functions	n.a.			5
number zero basis	n.a.			1
weight scale	n.a.	10		n.a.

Table 3: Hyperparameters for the Meta-World experiments.

	PPO	NDP	BBRL-PPO	BBRL-TRPL
number samples	16000			16
GAE λ	0.95			n.a.
discount factor	0.99			n.a.
ϵ_μ		n.a.		0.005
ϵ_Σ		n.a.		0.0005
optimizer			adam	
epochs	10			100
learning rate			3e-4	
use critic	True			False
epochs critic	10			n.a.
learning rate critic	3e-4			n.a.
number minibatches	32			n.a.
trust region loss weight		n.a.		10.0
entropy loss penalty			0	
normalized observations	True			False
normalized rewards	True			False
observation clip	10.0			n.a.
reward clip	10.0			n.a.
critic clip		0.2		n.a.
importance ratio clip		0.2		n.a.
hidden layers	[128, 128]			[32, 32]
hidden layers critic	[128, 128]			n.a.
hidden activation		tanh		relu
initial std	1.0			10.0
number basis functions	n.a.		5	
number zero basis	n.a.		1	
weight scale	n.a.	100		n.a.

Table 4: Hyperparameters for the hopper jumping experiments.

	PPO	BBRL-PPO	BBRL-TRPL
number samples	16384		320
GAE λ	0.95		n.a.
discount factor	0.99		n.a.
ϵ_μ		n.a.	0.005
ϵ_Σ		n.a.	0.0005
optimizer		adam	
epochs	10		100
learning rate	3e-4	1e-4	5e-5
use critic	True		False
epochs critic	10		n.a.
learning rate critic	3e-4		n.a.
number minibatches	32		n.a.
trust region loss weight		n.a.	25.0
entropy loss penalty		0	
normalized observations	True		False
normalized rewards	True		False
observation clip	10.0		n.a.
reward clip	10.0		n.a.
critic clip		0.2	n.a.
importance ratio clip		0.2	n.a.
hidden layers	[128, 128]		[32, 32]
hidden layers critic	[128, 128]		n.a.
hidden activation		tanh	
initial std	1.0		1.0
number basis functions	n.a.		5
number zero basis	n.a.		1

Table 5: Hyperparameters for the Beer Pong experiments.

	PPO	BBRL-PPO	BBRL-TRPL
number samples	16384		160
GAE λ	0.95		n.a.
discount factor	0.99		n.a.
ϵ_μ		n.a.	0.005
ϵ_Σ		n.a.	0.0005
optimizer		adam	
epochs	10		100
learning rate	1e-4	1e-4	5e-5
use critic	True		False
epochs critic	10		n.a.
learning rate critic	3e-4		n.a.
number minibatches	32		n.a.
trust region loss weight		n.a.	25.0
entropy loss penalty		0	
normalized observations	True		False
normalized rewards	True		False
observation clip	10.0		n.a.
reward clip	10.0		n.a.
critic clip		0.2	n.a.
importance ratio clip		0.2	n.a.
hidden layers	[128, 128]		[32, 32]
hidden layers critic	[128, 128]		n.a.
hidden activation		tanh	
initial std	1.0		1.0
number basis functions	n.a.		2
number zero basis	n.a.		2

Table 6: Hyperparameters for the Table Tennis experiments.

	PPO	BBRL-PPO	BBRL-TRPL
number samples	16000		200
GAE λ		0.95	n.a.
discount factor		0.99	n.a.
ϵ_μ		n.a.	0.0005
ϵ_Σ		n.a.	0.00005
optimizer		adam	
epochs	10		100
learning rate	3e-4	1e-4	3e-4
use critic		True	
epochs critic	10		100
learning rate critic	3e-4	1e-4	3e-4
number minibatches	32		n.a.
trust region loss weight		n.a.	25.0
entropy loss penalty		0	
normalized observations	True		False
normalized rewards	True		False
observation clip	10.0		n.a.
reward clip	10.0		n.a.
critic clip		0.2	n.a.
importance ratio clip		0.2	n.a.
hidden layers	[256, 256]		[256]
hidden layers critic	[256, 256]		[256]
hidden activation		tanh	
initial std		1.0	
number basis functions	n.a.		3
number zero basis	n.a.		1