

APPENDIX

A ALGORITHM

Algorithm 1 DrQ-v2: Improved data-augmented RL.**Inputs:** $f_\xi, \pi_\phi, Q_{\theta_1}, Q_{\theta_2}$: parametric networks for encoder, policy, and Q-functions respectively.

aug: random shifts image augmentation.

 $\sigma(t)$: scheduled standard deviation for the exploration noise defined in Equation (3). T, B, α, τ, c : training steps, mini-batch size, learning rate, target update rate, clip value.**Training routine:****for** each timestep $t = 1..T$ **do** $\sigma_t \leftarrow \sigma(t)$

▷ Compute stddev for the exploration noise

 $\mathbf{a}_t \leftarrow \pi_\phi(f_\xi(\mathbf{x}_t)) + \epsilon$ and $\epsilon \sim \mathcal{N}(0, \sigma_t^2)$

▷ Add noise to the deterministic action

 $\mathbf{x}_{t+1} \sim P(\cdot | \mathbf{x}_t, \mathbf{a}_t)$

▷ Run transition function for one step

 $\mathcal{D} \leftarrow \mathcal{D} \cup (\mathbf{x}_t, \mathbf{a}_t, R(\mathbf{x}_t, \mathbf{a}_t), \mathbf{x}_{t+1})$

▷ Add a transition to the replay buffer

UPDATECRITIC(\mathcal{D}, σ_t)UPDATEACTOR(\mathcal{D}, σ_t)**end for****procedure** UPDATECRITIC(\mathcal{D}, σ) $\{(\mathbf{x}_t, \mathbf{a}_t, r_{t:t+n-1}, \mathbf{x}_{t+n})\} \sim \mathcal{D}$ ▷ Sample a mini batch of B transitions $\mathbf{h}_t, \mathbf{h}_{t+n} \leftarrow f_\xi(\text{aug}(\mathbf{x}_t)), f_\xi(\text{aug}(\mathbf{x}_{t+n}))$

▷ Apply data augmentation and encode

 $\mathbf{a}_{t+n} \leftarrow \pi_\phi(\mathbf{h}_{t+n}) + \epsilon$ and $\epsilon \sim \text{clip}(\mathcal{N}(0, \sigma^2))$

▷ Sample action

Compute $\mathcal{L}_{\theta_1, \xi}$ and $\mathcal{L}_{\theta_2, \xi}$ using Equation (1)

▷ Compute critic losses

 $\xi \leftarrow \xi - \alpha \nabla_\xi (\mathcal{L}_{\theta_1, \xi} + \mathcal{L}_{\theta_2, \xi})$

▷ Update encoder weights

 $\theta_k \leftarrow \theta_k - \alpha \nabla_{\theta_k} \mathcal{L}_{\theta_k, \xi} \quad \forall k \in \{1, 2\}$

▷ Update critic weights

 $\bar{\theta}_k \leftarrow (1 - \tau)\bar{\theta}_k + \tau\theta_k \quad \forall k \in \{1, 2\}$

▷ Update critic target weights

end procedure**procedure** UPDATEACTOR(\mathcal{D}, σ) $\{(\mathbf{x}_t)\} \sim \mathcal{D}$ ▷ Sample a mini batch of B observations $\mathbf{h}_t \leftarrow f_\xi(\text{aug}(\mathbf{x}_t))$

▷ Apply data augmentation and encode

 $\mathbf{a}_t \leftarrow \pi_\phi(\mathbf{h}_t) + \epsilon$ and $\epsilon \sim \text{clip}(\mathcal{N}(0, \sigma^2))$

▷ Sample action

Compute \mathcal{L}_ϕ using Equation (2)

▷ Compute actor loss

 $\phi \leftarrow \phi - \alpha \nabla_\phi \mathcal{L}_\phi$

▷ Update actor's weights only

end procedure

B BENCHMARKS

We classify a set of 24 continuous control tasks from DMC (Tassa et al., 2018) into *easy*, *medium*, and *hard* benchmarks and provide a summary for each task in Table 1.

Task	Traits	Difficulty	Allowed Steps	$\dim(\mathcal{S})$	$\dim(\mathcal{A})$
Cartpole Balance	balance, dense	easy	1×10^6	4	1
Cartpole Balance Sparse	balance, sparse	easy	1×10^6	4	1
Cartpole Swingup	swing	dense	1×10^6	4	1
Cup Catch	swing, catch, sparse	easy	1×10^6	8	2
Finger Spin	rotate, dense	easy	1×10^6	6	2
Hopper Stand	stand, dense	easy	1×10^6	14	4
Pendulum Swingup	swing, sparse	easy	1×10^6	2	1
Walker Stand	stand, dense	easy	1×10^6	18	6
Walker Walk	walk, dense	easy	1×10^6	18	6
Acrobot Swingup	diff. balance, dense	medium	3×10^6	4	1
Cartpole Swingup Sparse	swing, sparse	medium	3×10^6	4	1
Cheetah Run	run, dense	medium	3×10^6	18	6
Finger Turn Easy	turn, sparse	medium	3×10^6	6	2
Finger Turn Hard	turn, sparse	medium	3×10^6	6	2
Hopper Hop	move, dense	medium	3×10^6	14	4
Quadruped Run	run, dense	medium	3×10^6	56	12
Quadruped Walk	walk, dense	medium	3×10^6	56	12
Reach Duplo	manipulation, sparse	medium	3×10^6	55	9
Reacher Easy	reach, dense	medium	3×10^6	4	2
Reacher Hard	reach, dense	medium	3×10^6	4	2
Walker Run	run, dense	medium	3×10^6	18	6
Humanoid Stand	stand, dense	hard	30×10^6	54	21
Humanoid Walk	walk, dense	hard	30×10^6	54	21
Humanoid Run	run, dense	hard	30×10^6	54	21

Table 1: A detailed description of each tasks in our *easy*, *medium*, and *hard* benchmarks.

C FULL COMPARISON TO MODEL-FREE METHODS

C.1 FULL RESULTS ON EASY BENCHMARK

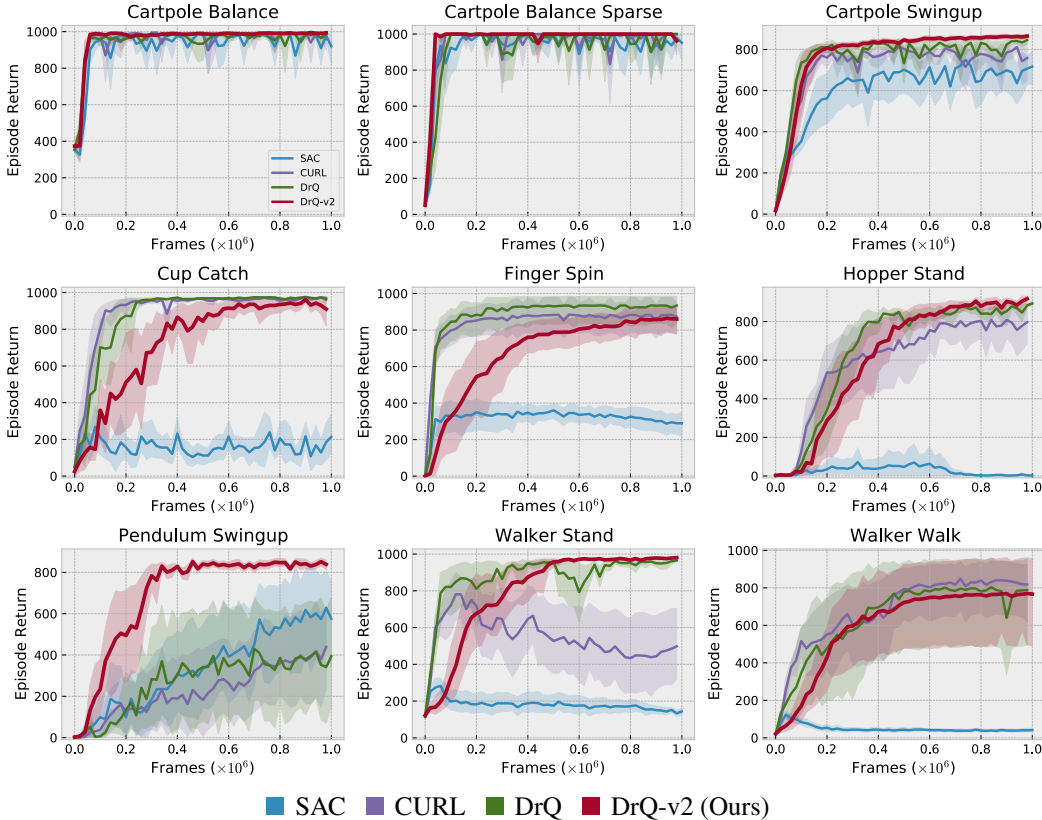


Figure 6: The *easy* benchmark consists of 9 tasks, where performance gains have been largely saturated by prior work. Still, DrQ-v2 is able to match sample complexity of the baselines. We note, that evaluation on these tasks is done for completeness reasons only, and encourage RL practitioners to refrain from using them for benchmarking purposes in future research.

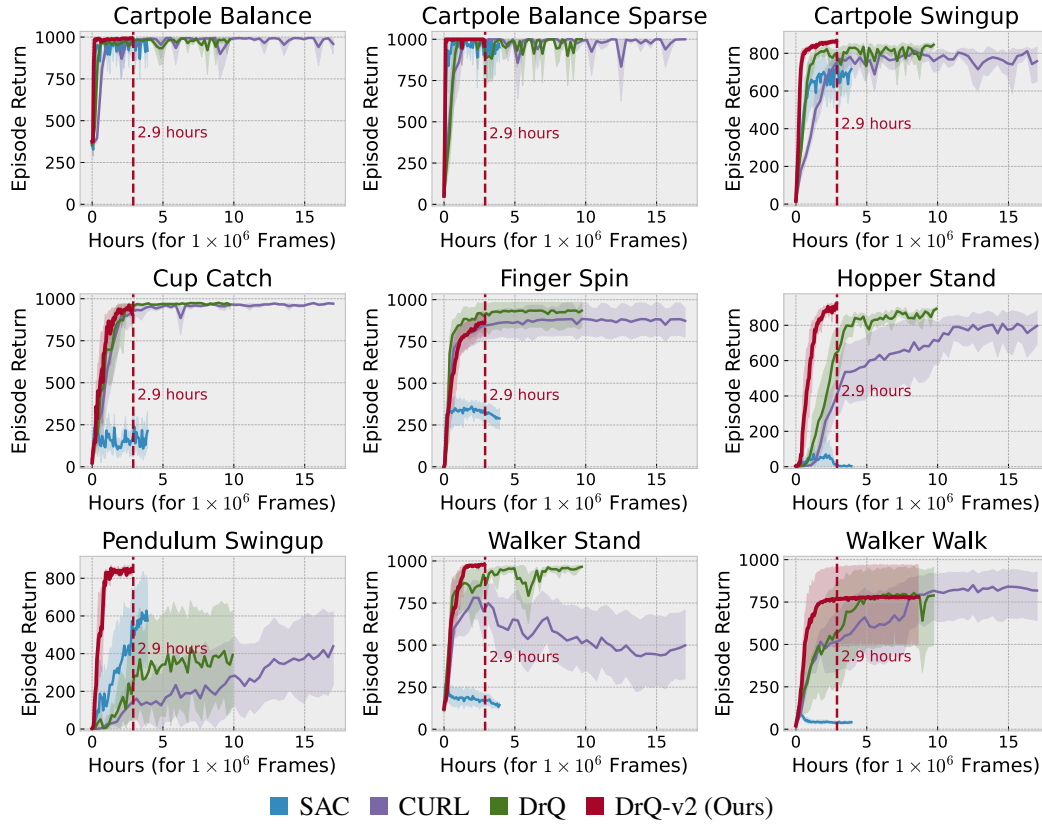


Figure 7: DrQ-v2 requires less than 3 hours to solve tasks from the *easy* benchmark, reaching a throughput of 96 FPS on a single NVIDIA V100 GPU.

C.2 FULL RESULTS ON MEDIUM BENCHMARK

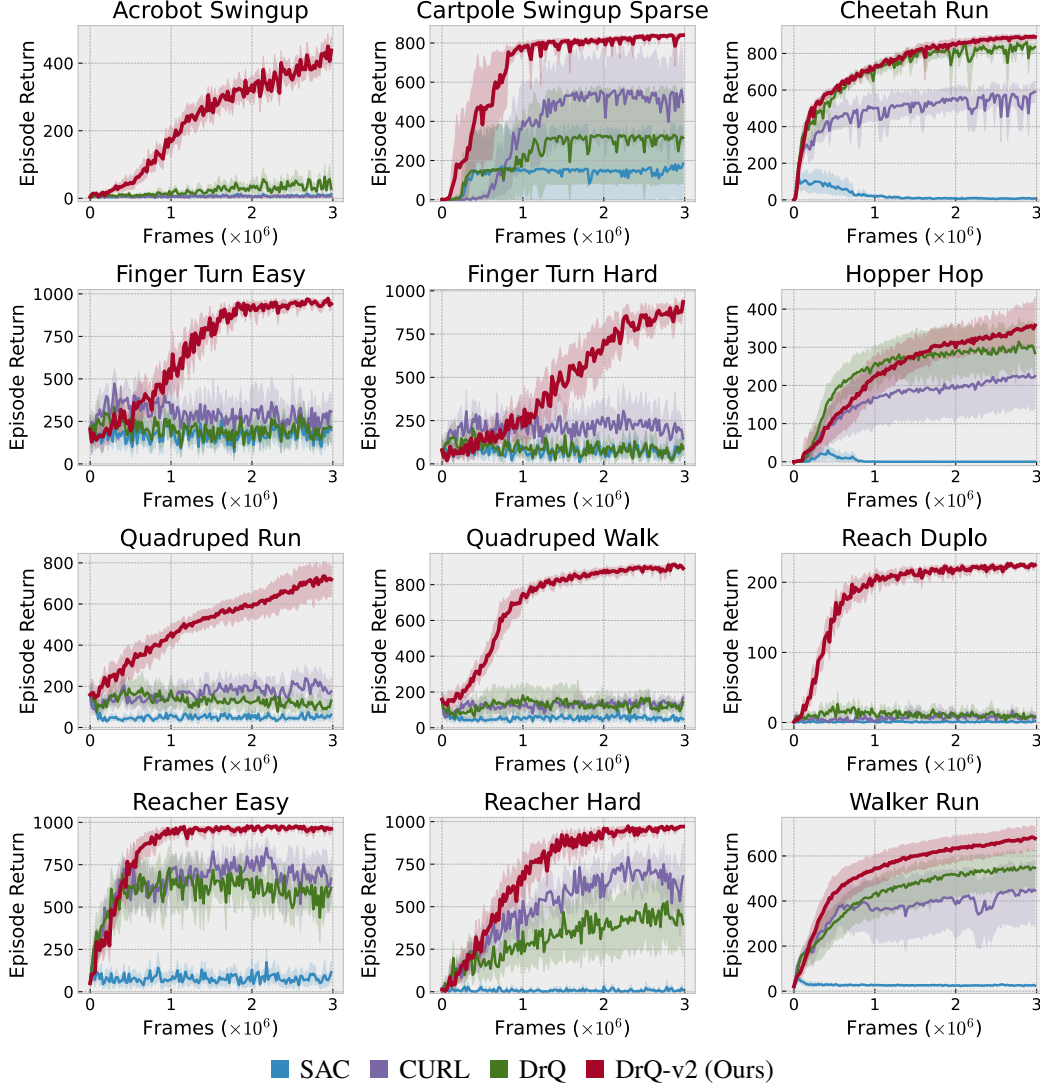


Figure 8: The *medium* benchmark consists of 12 complex control tasks that offer various challenges, including complex dynamics, sparse rewards, hard exploration, and more. DrQ-v2 demonstrates favorable sample efficiency and comfortably outperforms leading model-free baselines.

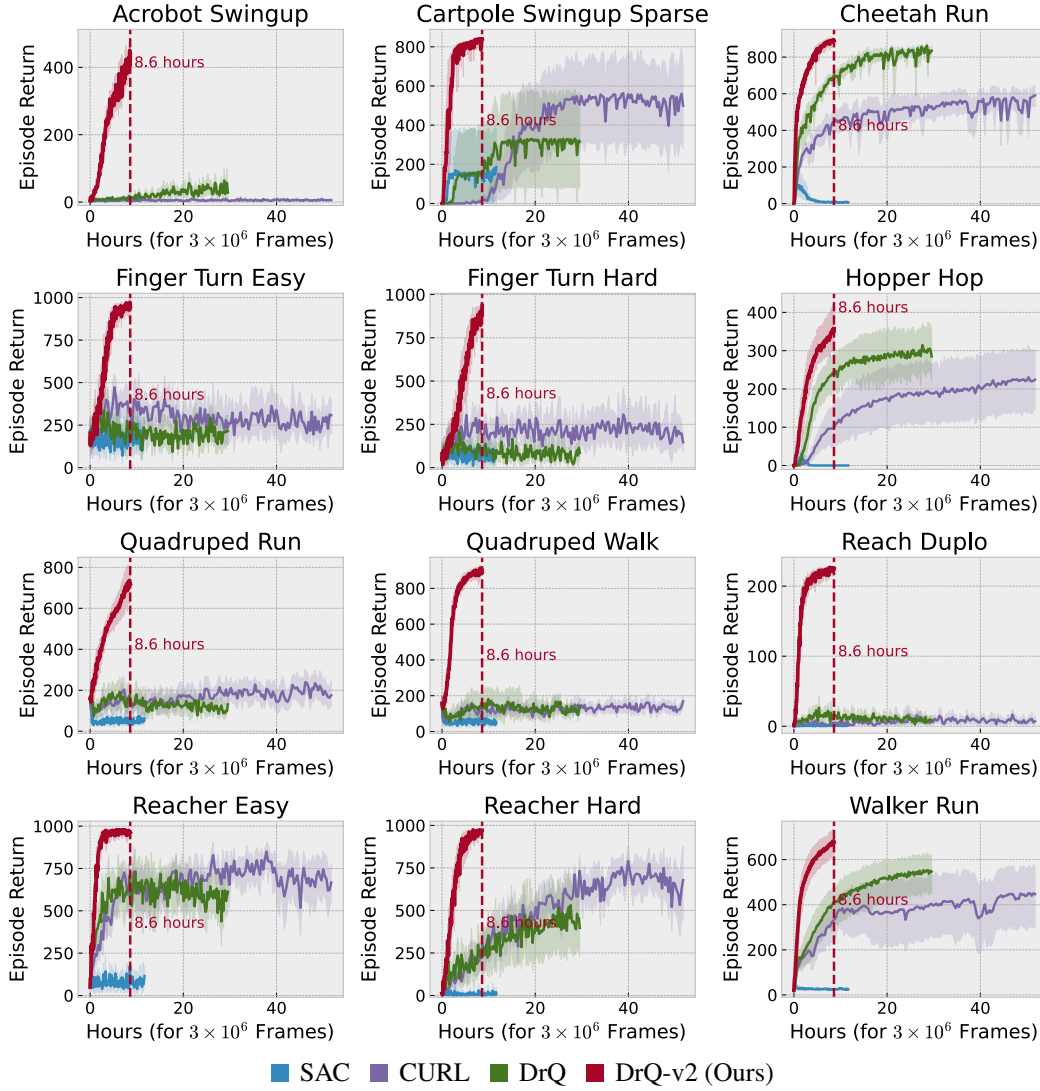


Figure 9: DrQ-v2 requires less than 9 hours to solve tasks from the *medium* benchmark, reaching a throughput of 96 FPS on a single NVIDIA V100 GPU.

C.3 FULL RESULTS ON HARD BENCHMARK

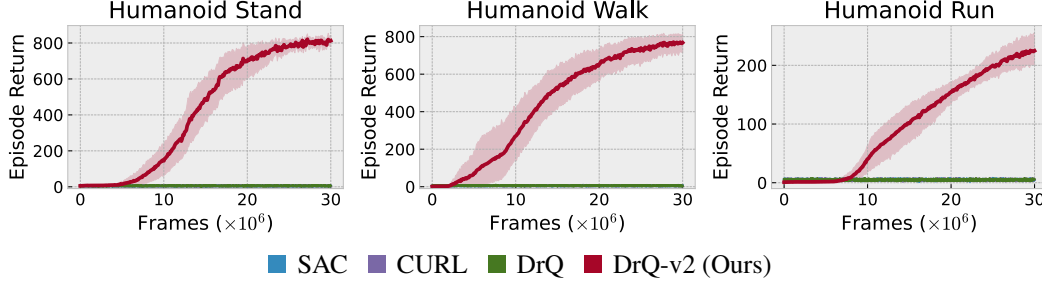


Figure 10: The *hard* benchmark consists of three humanoid locomotion tasks: *stand*, *walk*, and *run*. These three represent particularly hard exploration challenges, being previously unsolvable by model-free methods. The training speed of DrQ-v2 was key to solving the task, since it allowed for extensive investigation of different variations, resulting in the discovery of effective strategies.

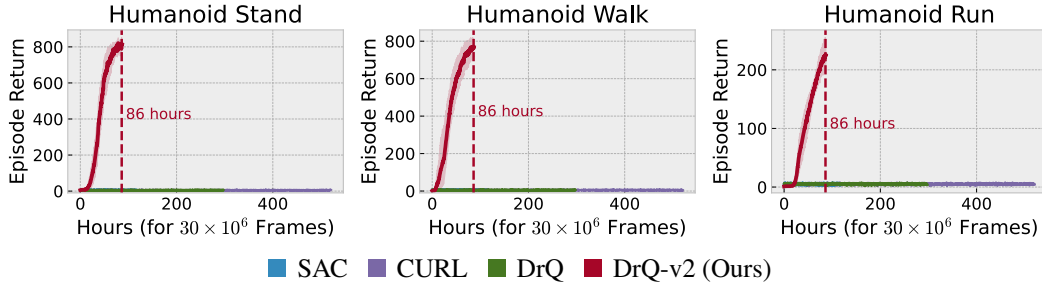


Figure 11: DrQ-v2 requires around 86 hours to solve tasks from the *hard* benchmark, reaching a throughput of 96 FPS on a single NVIDIA V100 GPU.

D FULL COMPARISON TO MODEL-BASED METHODS

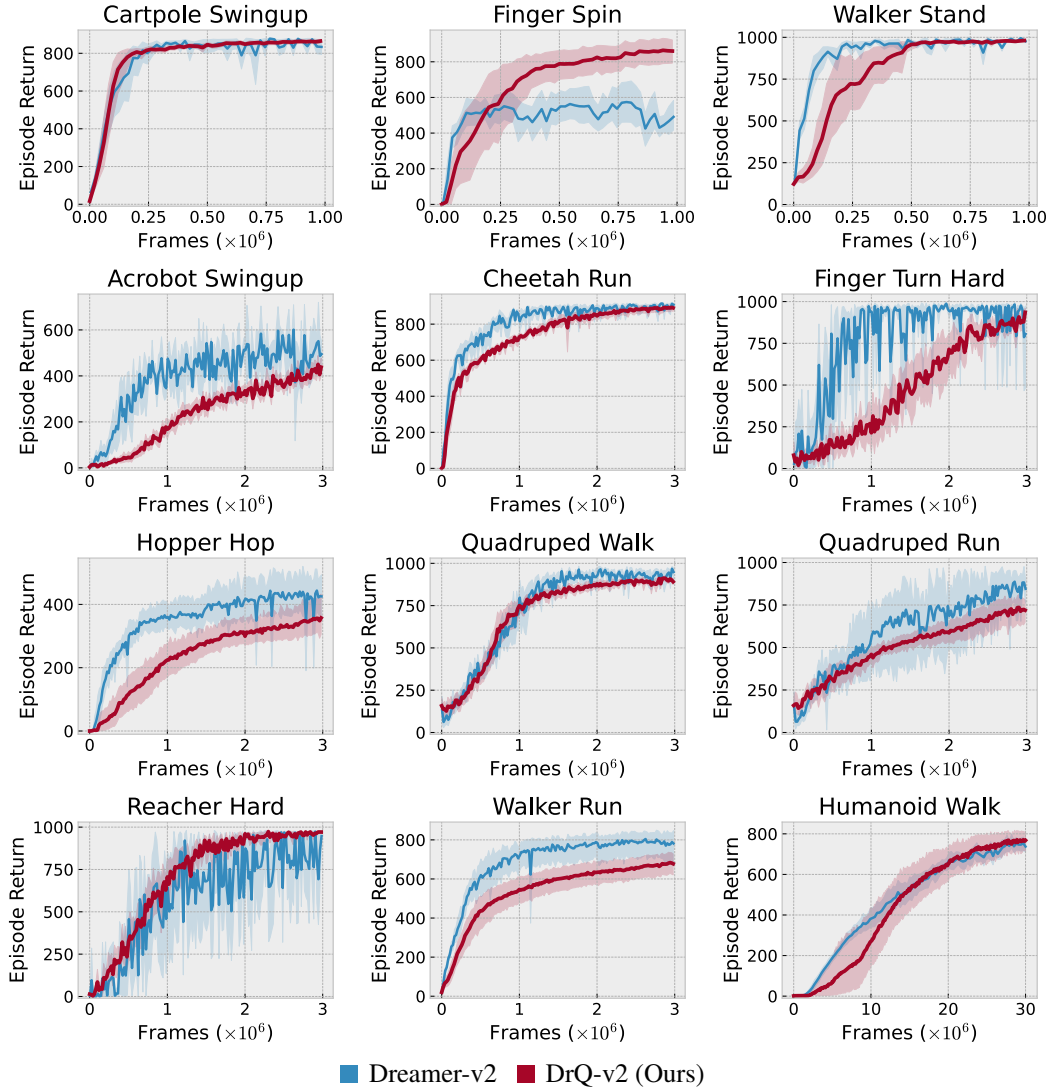


Figure 12: In many cases model-free DrQ-v2 can rival model-based Dreamer-v2 in sample efficiency. There are several tasks, however, where Dreamer-v2 performs better.

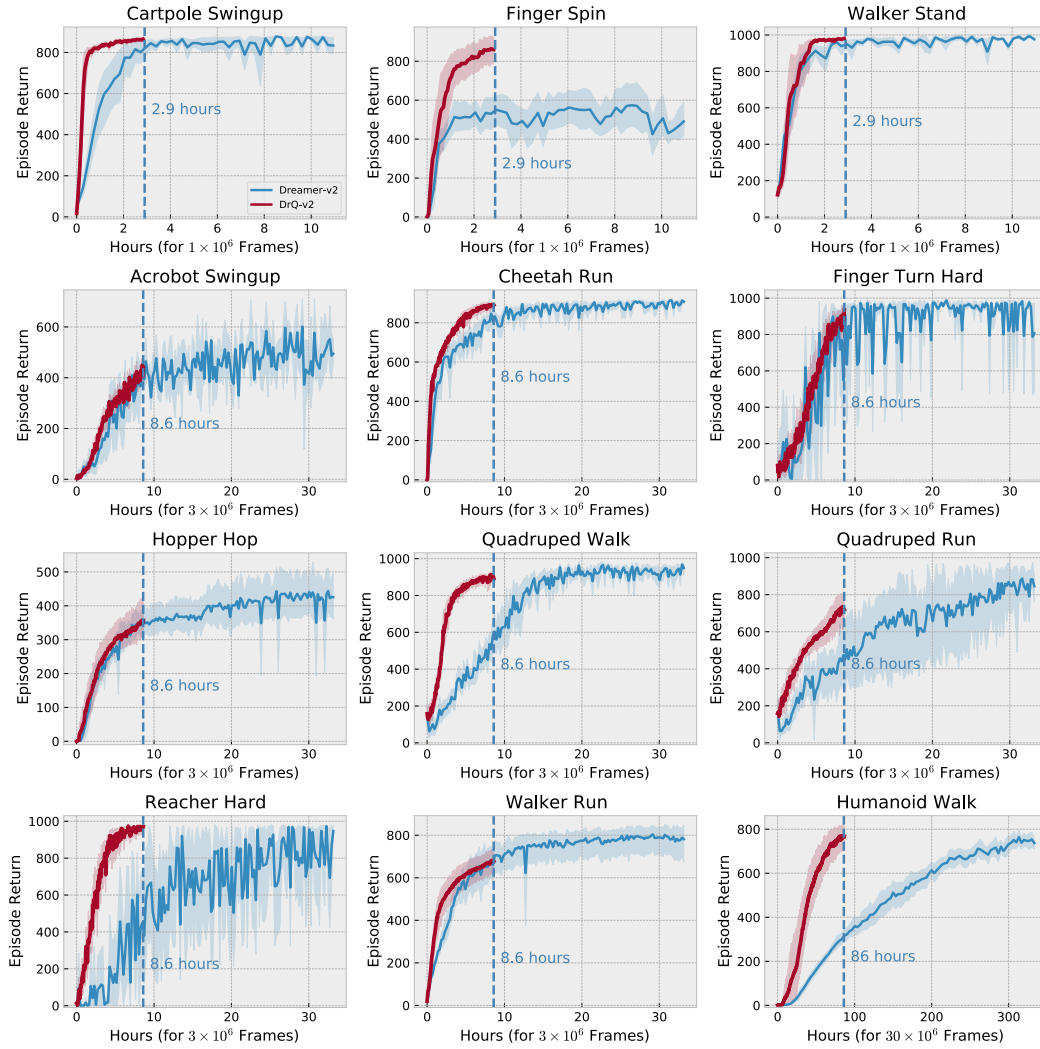


Figure 13: Model-based Dreamer-v2 performs more computations than model-free DrQ-v2. This allows DrQ-v2 to train faster in terms of wall-clock time and outperform Dreamer-v2 in this aspect.

E HYPER-PARAMETERS

The full list of hyper-parameters is presented in Table 2. While we tried to keep the settings identical for each of the task, there are a few specific deviations for some tasks.

Walker Stand/Walk/Run For all three tasks we use mini-batch size of 512 and n -step return of 1.

Cartpole Swingup Sparse stddev. schedule is set to 1.0 to facilitate stronger exploration in the sparse reward setting.

Quadruped Run We set the replay buffer size to 10^5 .

Humanoid Stand/Walk We set learning rate to 8×10^{-5} and increase features dim. to 100.

Table 2: A default set of hyper-parameters used in our experiments.

Parameter	Setting
Replay buffer capacity	10^6
Action repeat	2
Seed frames	4000
Exploration steps	2000
n -step returns	3
Mini-batch size	256
Discount γ	0.99
Optimizer	Adam
Learning rate	10^{-4}
Agent update frequency	2
Critic Q-function soft-update rate τ	0.01
Features dim.	50
Hidden dim.	1024
Exploration stddev. clip	0.3
Exploration stddev. schedule	easy: linear(1.0, 0.1, 100000) medium: linear(1.0, 0.1, 500000) hard: linear(1.0, 0.1, 2000000)