

## Appendix: High-Throughput Synchronous Deep RL

In this appendix we first provide the proofs for Claim 1 (Sec. A) and Claim 2 (Sec. B). We then discuss delayed gradient updates (Sec. C), additional ablation studies (Sec. D), comparison with additional baselines (Sec. E), implementation details (Sec. F), metrics (Sec. G) and provide all the training curves (Sec. H).

### A Proof of Claim 1

**Claim 1.** Consider collecting  $K$  states using  $n$  parallel environments. Let  $X_i^{(j)}$  denote the time for the  $j^{\text{th}}$  environment to perform its  $i^{\text{th}}$  step. Suppose  $X_i^{(j)}$  is independent and identically distributed (i.i.d.) and  $\sum_{i=1}^{\alpha} X_i^{(j)}$  follows a Gamma distribution with shape  $\alpha$  and rate  $\beta$ . Assume the computation time of each actor consistently takes time  $c$ . Given these assumptions, the expected time  $\mathbb{E}[T_{\text{total}}^{n,K}]$  to generate  $K$  states approaches

$$\frac{K}{n\alpha} \left( \frac{\gamma}{\beta} \left( 1 + \frac{\alpha-1}{\beta F^{-1}(1-\frac{1}{n})} \right) + F^{-1}(1-\frac{1}{n}) \right) + \frac{Kc}{n}, \quad (7)$$

where  $F^{-1}$  is the inverse cumulative distribution function (inverse CDF) of a gamma distribution with shape  $\alpha$  and rate  $\beta$ , i.e.,  $\text{Gamma}(\alpha, \beta)$ , and  $\gamma$  is the Euler-Mascheroni constant.

*Proof.* As the environments synchronize every  $\alpha$  steps, we need  $\frac{K}{n\alpha}$  synchronizations to finish the  $K$  steps. Let  $T_l$  denote the time required for the  $l^{\text{th}}$  synchronization. We have  $\mathbb{E}[T_{\text{total}}^{n,K}] = \sum_{l=1}^{\frac{K}{n\alpha}} \mathbb{E}[T_l]$ . Note  $\mathbb{E}[T_l] = \mathbb{E}[\max_j \sum_{i=1}^{\alpha} X_i^{(j)}] + \alpha c \forall l$ . By assumption we know that  $Y_j \triangleq \sum_{i=1}^{\alpha} X_i^{(j)}$  follows a gamma distribution with shape  $\alpha$  and rate  $\beta$ . By extreme value theory [5, 9], suppose  $X_l^{(j)} \sim \text{Gamma}(\alpha, \beta)$ , then  $\mathbb{E}[\max_j X_l^{(j)}] \simeq \frac{\gamma}{\beta} \left( 1 + \frac{\alpha-1}{\beta F^{-1}(1-\frac{1}{n})} \right) + F^{-1}(1-\frac{1}{n})$ , where  $F^{-1} = \inf\{x \in \mathcal{R} : F(x) \geq q\}$ .  $F(x)$  is the CDF of  $\text{Gamma}(\alpha, \beta)$ , and  $\gamma$  is the Euler-Mascheroni constant. By plugging the obtained approximation into  $\sum_{l=1}^{\frac{K}{n\alpha}} \mathbb{E}[T_l]$ , the result follows.  $\square$

In Claim 1, we assume the sum of steptimes (synchronization time) to follow a Gamma distribution. We empirically verify this assumption. In Fig. A1, we show the histogram of synchronization time (sum of every 100 step times) on ‘3 vs. 1 w/ keeper’. Furthermore, we perform a Kolmogorov-Smirnov goodness-of-fit test, with a significance-level of 0.05 and D-statistics of 0.04. We find the empirical data is consistent with the assumed Gamma distribution.

### B Proof of Claim 2

**Claim 2.** Consider asynchronous parallel actor-learner systems, such as GA3C and IMPALA. Suppose the system has  $n$  actors and each actor is sending data to the data queue following an i.i.d. Poisson distribution with rate  $\lambda_0$ . The learners consume data following an exponential distribution

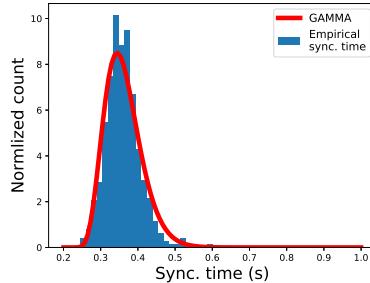


Figure A1: Empirical synchronization time.

with rate  $\mu$ . Let  $L$  denote the latency between the behavior policy and the target policy. Then, we have  $\mathbb{E}[L] = \frac{n\rho_0}{1-n\rho_0}$ , where utilization  $\rho_0 = \frac{\lambda_0}{\mu}$ .

*Proof.* Observe that the latency  $L$  is equal to the length of the data queue.  $n$  actors send data to the queue with rate  $n\lambda_0$  in total. Let  $P_i$  denote the probability that there are  $i$  data points in the queue. To be stable, the system must satisfy the balance equations

$$n\lambda_0 P_0 = \mu P_1 \quad (\text{A1})$$

$$(n\lambda_0 + \mu)P_j = n\lambda_0 P_{j-1} + \mu P_{j+1}, \quad j \geq 1. \quad (\text{A2})$$

Note Eq. (A1) and Eq. (A2) reduce to

$$(n\lambda_0)P_j = \mu P_{j+1}, j \geq 0, \quad (\text{A3})$$

or  $P_{j+1} = n\rho_0 P_j, j \geq 0$  from which we recursively obtain

$$P_j = (n\rho_0)^j P_0. \quad (\text{A4})$$

Using the fact that  $1 = \sum_{j=0}^{\infty} P_j = P_0 \sum_{j=0}^{\infty} (n\rho_0)^j$ , we observed that there is a solution if and only if  $n\rho_0 < 1$ , in which case  $1 = P_0(1 - n\rho_0)^{-1}$ , or

$$P_0 = 1 - n\rho_0. \quad (\text{A5})$$

Therefore, we have

$$P_j = (n\rho_0)^j (1 - n\rho_0), \quad (\text{A6})$$

which follows a geometric distribution with success probability  $(1 - n\rho_0)$ . Therefore, we have  $\mathbb{E}[L] = \frac{n\rho_0}{1-n\rho_0}$  which concludes the proof.  $\square$

## C Delayed Gradient

A delayed stochastic gradient descent performs the following update:  $\theta_t = \theta_{t-1} - \alpha_t \nabla \ell(x_{t-\tau}; \theta_{t-\tau})$ . The algorithm is identical to the standard stochastic gradient descent, except that gradients are delayed by a time step of  $\tau$ .

Consider a loss function of the form

$$\mathcal{L}(\theta) \triangleq \sum_{t=1}^T \ell(x_t; \theta). \quad (\text{A7})$$

We are interested in analyzing the convergence rate of  $\theta$  to the optimal parameters  $\theta^* \triangleq \arg \min_{\theta} \mathcal{L}(x_t; \theta)$ . Following Langford et al. [17], we assume the following: (a)  $\ell$  is convex, (b)  $L$ -Lipschitz, i.e.,  $\|\nabla \ell(x, \theta)\| \leq L$ , (c)  $x_t$  is drawn i.i.d. following a uniform distribution from a finite set  $X$ , (d)  $\max_{x, x' \in X} \frac{1}{2} \|x - x'\|^2 \leq F^2$ , where  $F$  is a constant, and (e) the learning rate of the delayed stochastic gradient descent is  $\frac{\sigma}{\sqrt{t-\tau}}$ , where  $\sigma^2 = \frac{F^2}{2\tau L^2}$ , then

$$\sum_{t=1}^T \ell(x_t, \theta_t) - \ell(x_t, \theta^*) \leq 4FL\sqrt{\tau T}. \quad (\text{A8})$$

Dividing both sides by  $T$ , we have

$$\frac{1}{T} \sum_{t=1}^T \ell(x_t, \theta_t) - \ell(x_t, \theta^*) \leq 4FL\sqrt{\frac{\tau}{T}}. \quad (\text{A9})$$

Stated differently, the convergence rate is  $O(\sqrt{\frac{\tau}{T}})$ . For HTS-RL with on-policy RL algorithms, the delay is guaranteed to be one, i.e.,  $\tau = 1$ . Therefore, the convergence rate is  $O(\sqrt{\frac{1}{T}})$ . Note that in practice the aforementioned assumptions are typically not met due to the use of deep nets.

## D Ablation Study

	Our Delayed Gradient	Truncated I.S.	No Correction
BankHeist	<b>987</b>	881	877
Breakout	<b>415</b>	390	402
Seaquest	<b>1831</b>	1827	1784

Table A1: Average episode rewards of our delayed-gradient, truncated importance sampling, and no correction on Atari games.

### D.1 Delayed Gradient

In addition to the convergence rate bound of delayed gradient, we verify the effectiveness of delayed gradient empirically. We run HTS-RL with (1) delayed gradient, (2) truncated importance sampling, and (3) no correction on multiple Atari games. The results are summarized in Tab. A1, where the average rewards of 100 evaluation episodes are reported. Compared with truncated importance sampling and no correction, the one-step delayed gradient in HTS-RL achieves a higher reward, which underlines the suitability of the delayed gradient strategy.

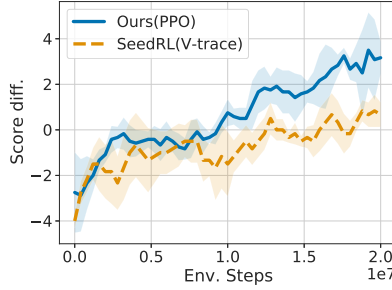


Figure A2: HTS-RL and SeedRL on GFootball 11 vs. 11 easy task.

Method	Kostrikov [14]	OpenAI Baselines [6]	rlpyt [29]	Ours
BankHeist	1382 $\pm$ 6	991 $\pm$ 14	1737 $\pm$ 39	<b>2111 <math>\pm</math> 21</b>
Beam Rider	1663 $\pm$ 14	1081 $\pm$ 18	2086 $\pm$ 32	<b>2586 <math>\pm</math> 14</b>
Breakout	1225 $\pm$ 12	829 $\pm$ 31	1508 $\pm$ 60	<b>1885 <math>\pm</math> 15</b>
Frostbite	1337 $\pm$ 8	962 $\pm$ 15	1803 $\pm$ 17	<b>1973 <math>\pm</math> 24</b>
Jamesbond	1353 $\pm$ 5	1014 $\pm$ 1	1991 $\pm$ 24	<b>2139 <math>\pm</math> 31</b>
Krull	1443 $\pm$ 6	1057 $\pm$ 11	2001 $\pm$ 29	<b>2657 <math>\pm</math> 16</b>
KFMaster	1532 $\pm$ 15	1056 $\pm$ 8	1979 $\pm$ 55	<b>2483 <math>\pm</math> 15</b>
MsPacman	1574 $\pm$ 9	1052 $\pm$ 3	1972 $\pm$ 13	<b>2364 <math>\pm</math> 5</b>
Qbert	1232 $\pm$ 13	953 $\pm$ 7	1621 $\pm$ 43	<b>1860 <math>\pm</math> 6</b>
Seaquest	1593 $\pm$ 10	946 $\pm$ 21	1918 $\pm$ 25	<b>2633 <math>\pm</math> 32</b>
S. Invader	1514 $\pm$ 20	1010 $\pm$ 7	1899 $\pm$ 32	<b>2318 <math>\pm</math> 12</b>
Star Gunner	1622 $\pm$ 19	1110 $\pm$ 5	2066 $\pm$ 24	<b>2616 <math>\pm</math> 25</b>

Table A2: SPS of different implementations of A2C.

## E Baselines

### E.1 A2C

To ensure the A2C implementation [14] we use is a strong baseline, we compare the speed of different versions of A2C, including our HTS-RL, Kostrikov [14], OpenAI baselines [6], and rlpyt [29], on Atari games. For a fair comparison, all methods use 16 parallel environment processes for data collection, and one GPU for model training/forwarding. For rlpyt, we use the most efficient ‘parallel-

	IMPALA	A2C/Ours
Unroll length	20	5
Batch size	32	-
Discount factor	0.99	0.99
Value loss coefficient	0.5	0.5
Entropy loss coefficient	0.01	0.01
RMSProp momentum	0.00	0.00
RMSProp $\epsilon$	0.01	0.00001
Learning rate	0.006	0.0007
Number of actors	16	4

Table A3: Hyper-parameters of IMPALA and A2C/Ours(A2C) in Atari experiments.

Method	IMPALA 16 actors (Baseline)	IMPALA 48 actors [16]
BankHeist	339	$\sim 300$
Beam Rider	4000	$\sim 4000$
Breakout	201	$\sim 130$
Frostbite	73	$\sim 70$
Jamesbond	82	$\sim 80$
Krull	2546	$\sim 2500$
KFMaster	9516	$\sim 8000$
MsPacman	807	$\sim 1300$
Qbert	4116	$\sim 4000$
Seaquest	458	$\sim 420$
S. Invader	1142	$\sim 2000$
Star Gunner	8560	$\sim 6000$

Table A4: Atari@20M environment steps. Since Küttler et al. [16] don’t report exact scores at 20M environment steps, we obtain their numbers from their plots and indicate that with a  $\sim$  symbol.

GPU’ mode. As shown in Tab. A2, Kostrikov’s A2C is a strong baseline, which achieves  $1.4\times$  higher SPS than OpenAI baselines. Also, HTS-RL consistently achieves higher SPS than rlpyt.

## E.2 SeedRL

SeedRL [8] is a recent work that reports results on GFootball ‘11 vs. 11 easy’ task. We compare HTS-RL with Seed RL (V-trace) [8] on Gfootball ‘11 vs. 11 easy.’ For a fair comparison, both HTS-RL and Seed RL use 16 parallel environment processes and one GPU. HTS-RL achieves 829 environment steps per second (SPS) while Seed RL achieves 609 SPS. After 20M steps of training, HTS-RL and Seed RL achieve a  $3.55 \pm 0.3$  and  $1.50 \pm 0.7$  score difference, respectively. The training curve is shown in Fig. A2.

# F Implementation Details

## F.1 Atari Game Experiments

In Atari experiments, we use the same neural network architectures as Espeholt et al. [7], Küttler et al. [16] for all three methods (IMPALA, A2C, Ours). The network has four hidden layers. The first layer is a convolutional layer with 32 filters of size  $8 \times 8$  and stride 4. The second layer is a convolutional layer with 64 filters of size  $4 \times 4$  and stride 2. The third layer is a convolutional layer with 64 filters of size  $3 \times 3$  and stride 1. The fourth layer is a fully connected layer with 512 hidden units. Following the hidden units are two sets of output. One provides a probability distribution over all valid actions. The other one provides the estimated value function. For ours (A2C) and A2C baseline, we use the same hyper-parameters as Kostrikov [14]. For IMPALA, we use the same hyper-parameters as Espeholt et al. [7], Küttler et al. [16]. We summarize the hyper-parameters in Tab. A3. Note Küttler et al. [16] deploy distributed IMLALA with 48 actors. However, in this work we target single machine parallel computing, and restrict ourselves to 16 parallel environments. For a fair comparison, we run all experiments with 16 parallel environments on a single machine. Importantly, while being downscaled to one machine, the reported IMPALA results match the results reported in the original paper [16]. Tab. A4 summarizes the results of our baseline and that reported by Küttler et al. [16].

## F.2 GFootball Experiments

In GFootball experiments, we use the CNN architecture of Kurach et al. [15] for all three methods (IMPALA, PPO, Ours). The network has four hidden layers. The first layer is a convolutional layer with 32 filters of size  $8 \times 8$  and stride 4. The second layer is a convolutional layer with 64 filters of

Method	IMPALA	IMPALA
	16 actors (Baseline)	500 actors [15]
Empty goal close	1.0	$\sim 0.99$
Empty goal	1.0	$\sim 0.85$
Run to score	0.80	$\sim 0.80$
RSK	0.05	$\sim 0.22$
PSK	0.20	$\sim 0.18$
RPSK	0.82	$\sim 0.41$
3 vs 1 w/ keeper	0.21	$\sim 0.20$
Corner	0.0	$\sim -0.1$
Counterattack easy	0.0	$\sim 0.0$
Counterattack hard	0.50	$\sim 0.0$
11 vs 11 w/ lazy Opp.	0.71	$\sim 0.38$

Table A5: GFootball Academy@5M environment steps. Since Kurach et al. [15] don’t report exact scores at 5M environment steps, we obtain their numbers from their plots and indicate that with a  $\sim$  symbol.

size  $4 \times 4$  and stride 2. The third layer is a convolutional layer with 64 filters of size  $3 \times 3$  and stride 1. The fourth layer is a fully connected layer with 512 hidden units. Following the hidden units are two sets of output. One provides a probability distribution over all valid actions. The other one provides the estimated value function. To be consistent with the official torch beast implementation [16], we use RMSProp for all methods. Regarding hyper-parameters, for ours (PPO) and PPO baseline, we mostly use the same hyper-parameters as Kurach et al. [15]. The only difference is that, instead of 512 steps, we unroll for 128 steps, which we found to give better results. For IMPALA, Kurach et al. [15] deploy distributed training with 500 actors. However, in this work we target single machine parallel computing, and restrict ourselves to 16 parallel environments. Therefore, for a fair comparison, we mostly follow the hyper-parameter settings of Kurach et al. [15], but decrease the number of actors and batch size. With only 16 actors and a smaller batch size, our baseline results match the results of IMPALA on GFootball environments reported by Kurach et al. [15]. Tab. A5 summarizes the results of our baseline and that reported by Kurach et al. [15]. The hyper-parameters are summarized in Tab. A6.

## G Final Time and Required Time Metrics

The results of all Atari experiments in *final time metric* and *required time metric* are summarized in Tab. A7 and Tab. A8. For *final time metric*, the time limit for each experiment is set to the time when IMPALA finishes training for 20M steps. For *required time metric*, we report the time to achieve average episode rewards of 40% and 80% of the A2C baseline episode rewards reported by Dhariwal et al. [6]. The results of all GFootball experiments in *final time metric* and *required time metric* are summarized in Tab. A9 and Tab. A10. For *final time metric*, the time limit for each experiment is set to the time when IMPALA finishes training for 5M steps. For *required time metric*, we report the time to achieve an average score of 0.4 and 0.8. As shown in Tab. A7 and Tab. A9, given the same amount of time, HTS-RL consistently achieves higher average rewards/scores than IMPALA and synchronous baselines. Moreover, as shown in Tab. A8 and Tab. A10, to achieve a target reward/score, HTS-RL consistently needs less time.

## H Training Curves

The training curves of all Atari and GFootball experiments in terms of time and number of environment steps are shown in Fig. A3, Fig. A4, Fig. A5, and Fig. A6. As shown in Fig. A4 and Fig. A6, HTS-RL does not trade data efficiency for higher throughput. While achieving much higher throughput, HTS-RL still maintains a data efficiency similar to synchronous baselines. As a result, HTS-RL consistently achieves higher rewards in shorter time than IMPALA and synchronous baselines across different environments (Fig. A3, Fig. A5).

	IMPALA	PPO / PPO(Ours)
Unroll length	32	128
Batch size	8	16
Discount factor	0.993	0.993
Value loss coefficient	0.5	0.5
Entropy loss coefficient	0.00087453	0.003
RMSProp momentum	0.00	0.00
RMSProp $\epsilon$	0.01	0.00001
Learning rate	0.00019896	0.000343
Number of actors	16	4

Table A6: Hyper-parameters of IMPALA and PPO/Ours(PPO) in GFootball experiments.

Method	IMPALA	A2C	Ours (A2C)
BankHeist	339 $\pm$ 10	775 $\pm$ 166	<b>942 <math>\pm</math> 100</b>
Beam Rider	4000 $\pm$ 690	4392 $\pm$ 134	<b>6995 <math>\pm</math> 420</b>
Breakout	201 $\pm$ 133	362 $\pm$ 29	<b>413 <math>\pm</math> 37</b>
Frostbite	73 $\pm$ 2	272 $\pm$ 14	<b>315 <math>\pm</math> 12</b>
Jamesbond	82 $\pm$ 10	438 $\pm$ 59	<b>474 <math>\pm</math> 88</b>
Krull	2546 $\pm$ 551	7560 $\pm$ 892	<b>7737 <math>\pm</math> 609</b>
KFMaster	9516 $\pm$ 3311	<b>30752 <math>\pm</math> 6641</b>	30020 $\pm$ 3559
MsPacman	807 $\pm$ 170	1236 $\pm$ 292	<b>1675 <math>\pm</math> 459</b>
Qbert	4116 $\pm$ 610	12479 $\pm$ 1965	<b>13682 <math>\pm</math> 1873</b>
Seaquest	458 $\pm$ 2	<b>1833 <math>\pm</math> 6</b>	1831 $\pm$ 7
S. Invader	<b>1142 <math>\pm</math> 207</b>	596 $\pm$ 69	731 $\pm$ 80
Star Gunner	8560 $\pm$ 918	41414 $\pm$ 3826	<b>52666 <math>\pm</math> 5182</b>

Table A7: Atari experiment in *final time metrics*: Average evaluation rewards achieved given limited training time.

Method (target reward 1 / target reward 2)	IMPALA	A2C	Ours (A2C)
BankHeist (480 / 960)	-/-	28.9/62.1	<b>18.9/116.8</b>
Beam Rider (1600 / 3200)	36.4/60.4	32.1/54.5	<b>10.3/30.9</b>
Breakout (160 / 320)	77.8/-	21.7/43.5	<b>17.7/38.9</b>
Frostbite (104 / 208)	-/-	5.0/10.0	<b>3.4/6.8</b>
Jamesbond (200 / 400)	-/-	39.4/49.2	<b>21.8/31.1</b>
Krull (3600 / 7200)	-/-	13.9/37.0	<b>7.5/37.6</b>
KFMaster (15200 / 30420)	-/-	39.5/192.6	<b>18.8/118.1</b>
MsPacman (880 / 1760)	75.8/-	49.3/160.3	<b>22.9/94.3</b>
Qbert (4000 / 8000)	94.2/-	53.3/83.8	<b>52.2/67.7</b>
Seaquest (640 / 1280)	-/-	6.7/28.4	<b>4.0/17.2</b>
Space. (240 / 480)	<b>9.65/19.9</b>	14.1/26.4	6.9/21.8
Star Gunner (8400 / 16800)	47.1/-	28.7/41.1	<b>17.8/25.4</b>

Table A8: Atari experiment in *required time metrics*: Required time (minutes) to achieve goal episode rewards (time required to achieve 40% rewards reported by Dhariwal et al. [6] / time required to achieve 80% rewards reported by Dhariwal et al. [6]). ‘-’ indicates that the method did not achieve the desired reward after 20M environment step training. Space.: Space invaders, KFMaster: KungFu Master.

Method	IMPALA	PPO	Ours (PPO)
Empty goal close	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>	<b>1.00 ± 0.00</b>
Empty goal	<b>1.00 ± 0.00</b>	0.89 ± 0.05	0.94 ± 0.04
Run to score	0.65 ± 0.42	0.89 ± 0.05	<b>0.93 ± 0.03</b>
RSK	0.03 ± 0.03	0.52 ± 0.21	<b>0.88 ± 0.06</b>
PSK	0.00 ± 0.00	0.05 ± 0.04	<b>0.41 ± 0.02</b>
RPSK	0.67 ± 0.06	0.49 ± 0.06	<b>0.80 ± 0.03</b>
3 vs 1 w/ keeper	0.23 ± 0.01	0.20 ± 0.09	<b>0.81 ± 0.02</b>
Corner	-0.10 ± 0.33	-0.06 ± 0.08	<b>0.03 ± 0.10</b>
Counterattack easy	0.00 ± 0.00	0.01 ± 0.01	<b>0.39 ± 0.02</b>
Counterattack hard	0.00 ± 0.00	0.01 ± 0.02	<b>0.53 ± 0.09</b>
11 vs 11 w/ lazy Opp.	0.46 ± 0.21	0.33 ± 0.07	<b>0.72 ± 0.09</b>

Table A9: GFootball experiments in *final time metrics*: Average evaluation scores achieved given limited training time. The time limit for each experiment is set to the time when IMPALA finishes training for 5M steps RSK: run to score w/ keeper, PSK: pass, shoot, w/ keeper, RPSK: run, pass, shoot, w/ keeper.

Method	IMPALA	PPO	Ours (PPO)
Empty goal close	1.7/2.6	5.4/15.5	<b>1.0/2.0</b>
Empty goal	8.4/11.7	12.8/19.2	<b>2.0/3.9</b>
Run to score	27.0/34.6	16.2/32.5	<b>6.3/11.4</b>
RSK	52.3/-	51.2/68.2	<b>11.5/18.8</b>
PSK	-/-	70.0/-	<b>38.8/-</b>
RPSK	22.3/ <b>25.4</b>	45.2/90.8	<b>13.5/27.1</b>
3 vs 1 w/ keeper	-/-	67.4/144.2	<b>15.9/25.6</b>
Corner	-/-	-/-	-/-
Counterattack easy	-/-	223.2/-	<b>91.3/-</b>
Counterattack hard	-/-	383.4/-	<b>61.8/-</b>
11 vs 11 w/ lazy Opp.	58.2/-	95.8/260.9	<b>14.4/72.1</b>

Table A10: GFootball experiments in *required time metrics*: required time (minutes) to achieve goal scores (time required to achieve score 0.4 / time required to achieve score 0.8). ‘-’ indicates that the method did not achieve the desired score after 5M environment step training. RSK: run to score w/ keeper, PSK: pass, shoot, w/ keeper, RPSK: run, pass, shoot, w/ keeper.

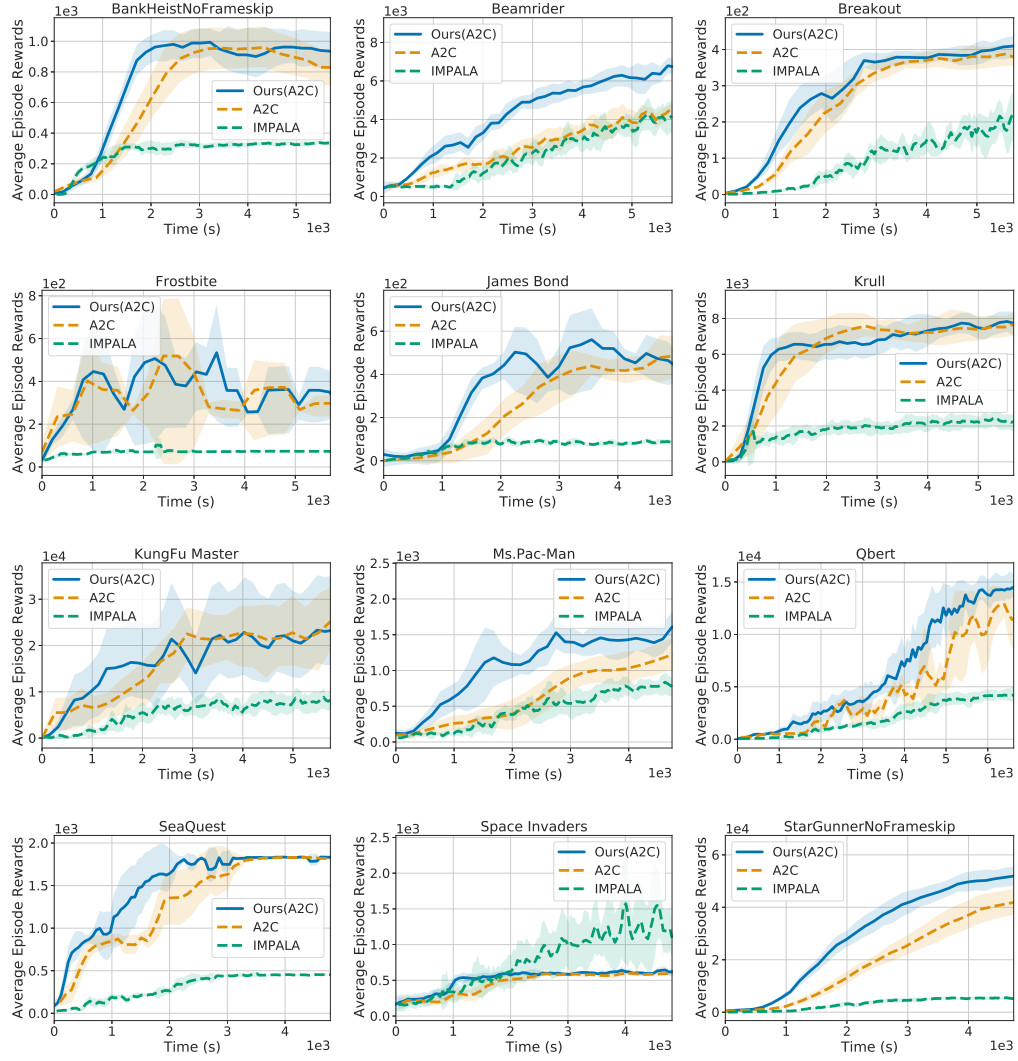


Figure A3: Atari: **Time** versus reward.



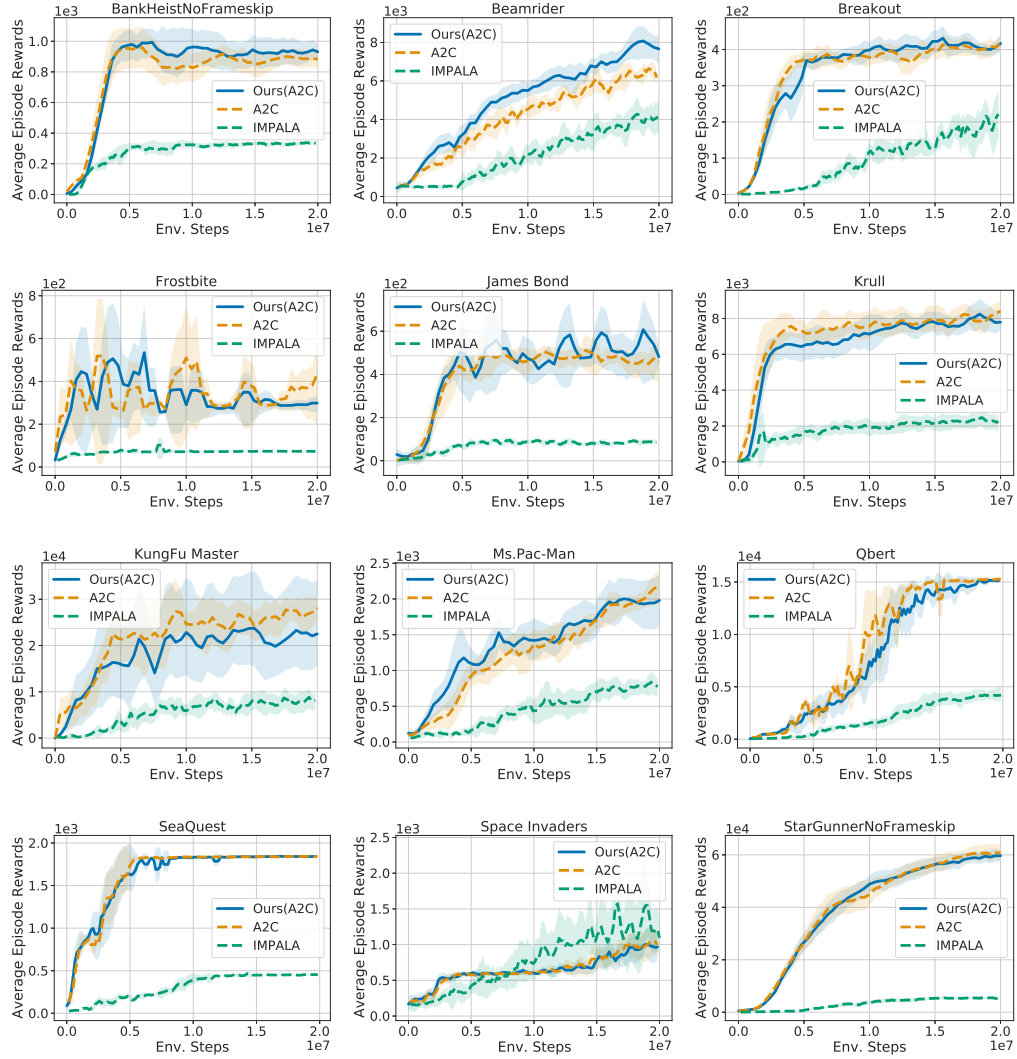


Figure A4: Atari: **Environment step** versus reward.

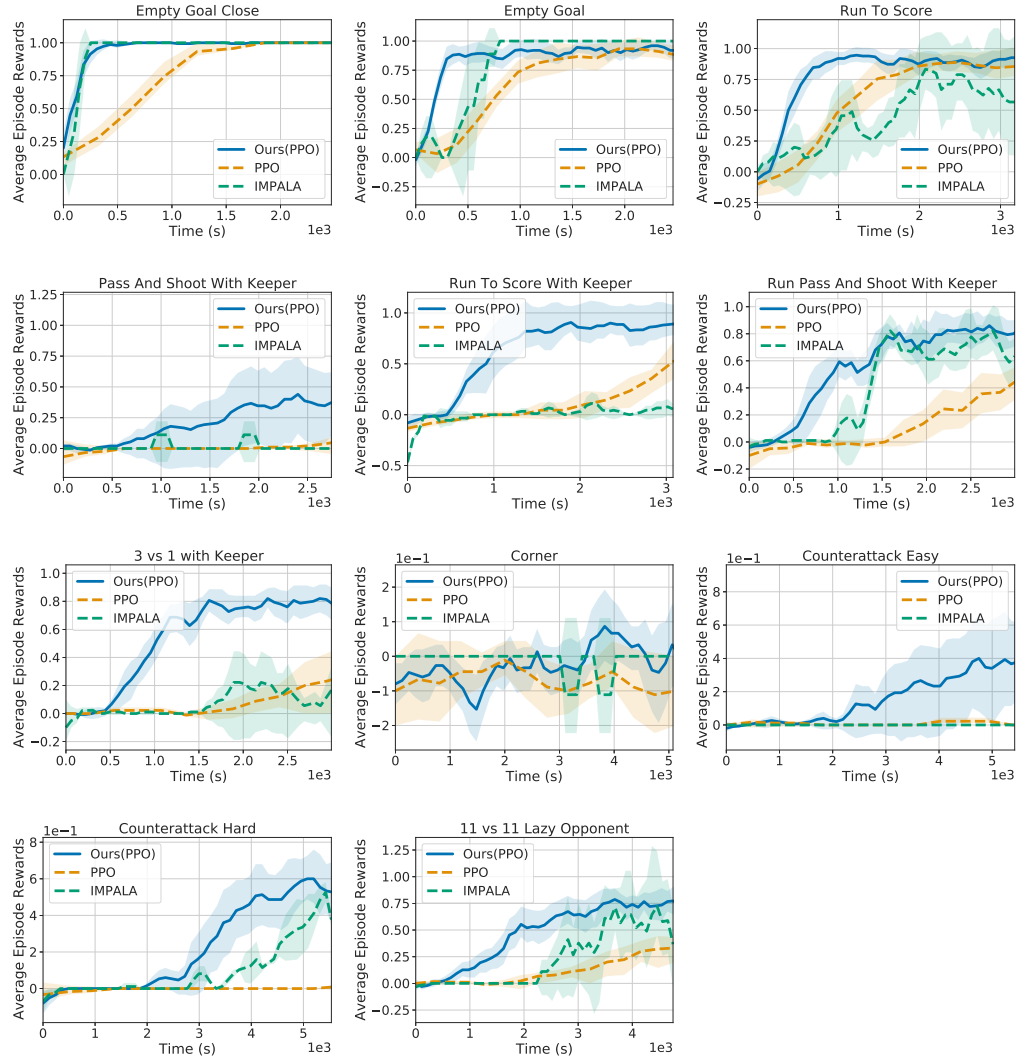


Figure A5: GFootball: **Time** versus reward.

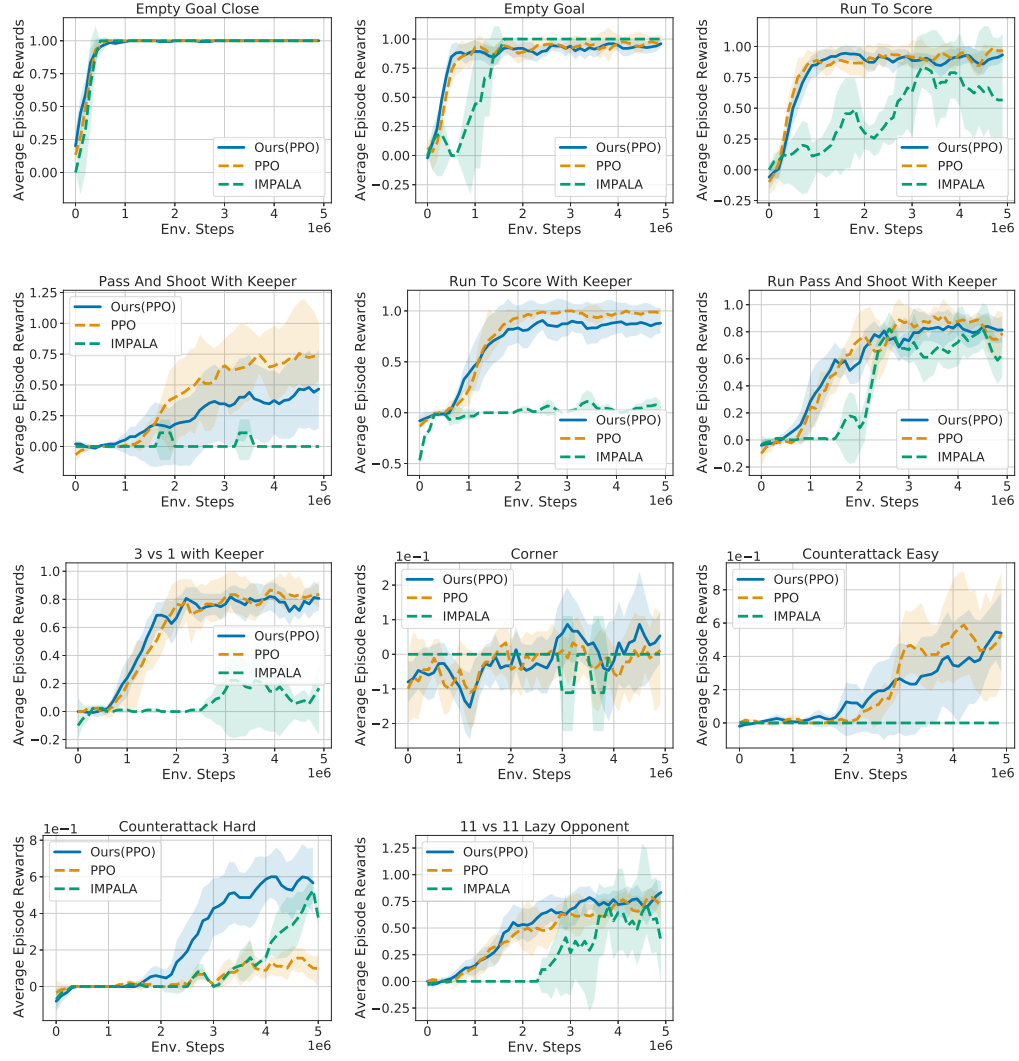


Figure A6: GFootball: **Environment step** versus reward.