# A. Appendix

## A.1. Algorithms

---

**Algorithm 1**

Online RL with decoupled ATC encoder (steps distinct from end-to-end RL in blue)

---

**Require:** $\theta_{ATC}, \phi_\pi$             $\triangleright$ ATC model parameters (encoder $f_\theta$ thru contrast $W$), policy parameters

  1: $\mathcal{S} \leftarrow \{\}$            $\triangleright$ replay buffer of observations

  2: $\bar{\theta}_{ATC} \leftarrow \theta_{ATC}$            $\triangleright$ initialize momentum encoder (conv and linear only)

  3: **repeat**

  4:      Sample environment and policy, through encoder:

  5:      **for** 1 to m **do**            $\triangleright$ a minibatch

  6:        $a \sim \pi(\cdot | f_\theta(s); \phi), s' \sim T(s,a), r \sim R(s,a,s')$

  7:        $\mathcal{S} \leftarrow \mathcal{S} \cup \{s\}$            $\triangleright$ store observations (delete oldest if full)

  8:        $s \leftarrow s'$

  9:      **end for**

10:      Update policy by given RL formula:            $\triangleright$ on- or off-policy

11:      **for** 1 to n **do**            $\triangleright$ given number RL updates per minibatch

12:        $\phi_\pi \leftarrow \phi_\pi + RL(s,a,s',r; \phi_\pi)$            $\triangleright$ stop gradient into encoder

13:      **end for**

14:      Update encoder (and contrastive model) by ATC:

15:      **for** 1 to p **do**

16:        $s, s_+ \sim \mathcal{S}$            $\triangleright$ sample observations: anchors and positives

17:        $\theta_{ATC} \leftarrow \theta_{ATC} - \lambda_{ATC} \nabla_{\theta_{ATC}} \mathcal{L}^{ATC}(s, s_+)$            $\triangleright$ ATC gradient update

18:        $\bar{\theta}_{ATC} \leftarrow (1-\tau)\bar{\theta}_{ATC} + \tau\theta_{ATC}$            $\triangleright$ update momentum encoder (conv and linear only)

19:      **end for**

20: **until** converged

21: **return** Encoder $f_\theta$ and policy $\pi_\phi$
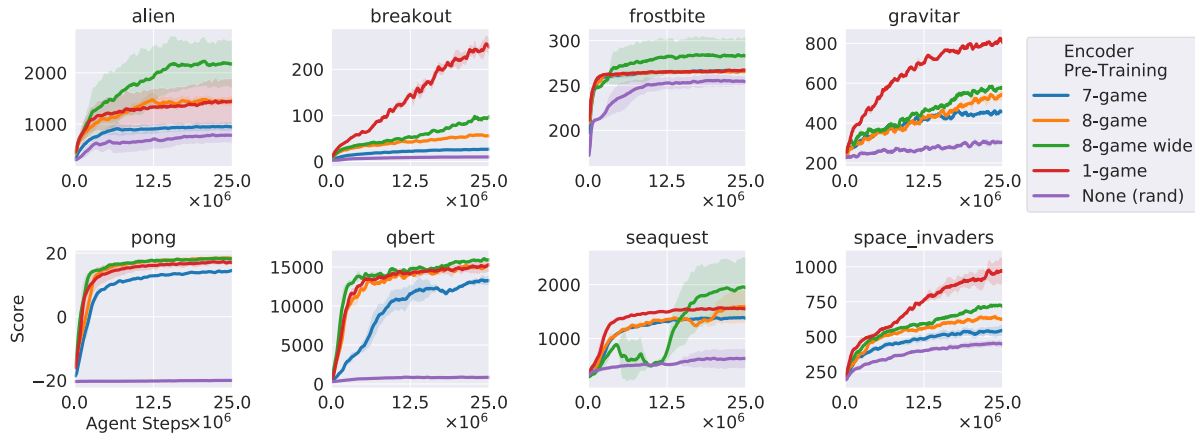
---

## A.2. Additional Figures



*Figure 13.* RL using multi-task encoders (all with weights frozen) for eight Atari games gives mixed performance, partially improved by increased network capacity (8-game-wide). Training on 7 games and testing on the held-out one yields diminished but non-zero performance, showing some limited feature transfer between games.
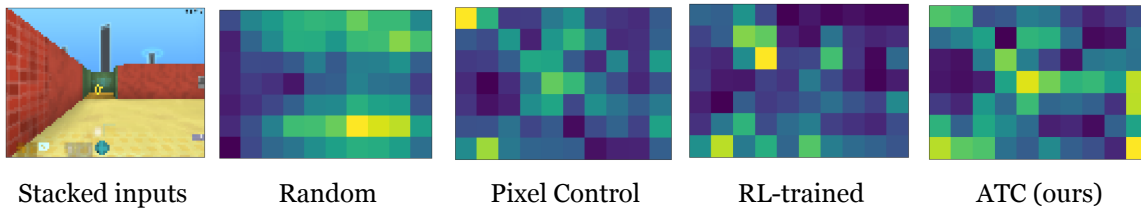


| Stacked inputs | Random | Pixel Control | RL-trained | ATC (ours) |

*Figure 14.* Attention map in LASERTAG. UL encoder with pixel control focuses on the score, while UL encoder with the proposed ATC focuses properly on the coin similar to RL-trained encoder.
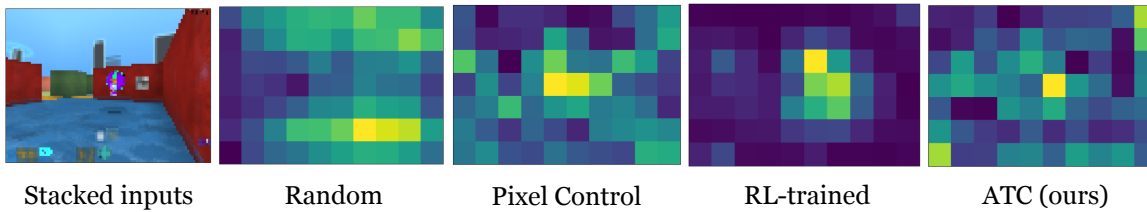


| Stacked inputs | Random | Pixel Control | RL-trained | ATC (ours) |

*Figure 15.* Attention map in the LASERTAG which shows that UL encoders focus properly on the enemy similar to RL-trained encoder.

## A.3. RL Settings

*Table 1.* DMControl, RAD-SAC Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| OBSERVATION RENDERING | $(84, 84)$, RGB |
| RANDOM SHIFT PAD | $\pm 4$ |
| REPLAY BUFFER SIZE | $1 \times 10^5$ |
| INITIAL STEPS | $1 \times 10^4$ |
| STACKED FRAMES | 3 |
| ACTION REPEAT | 2 (FINGER, WALKER) |
| | 8 (CARTPOLE) |
| | 4 (REST) |
| OPTIMIZER | ADAM |
| $(\beta_1, \beta_2) \to (f_\theta, \pi_\psi, Q_\phi)$ | $(.9, .999)$ |
| $(\beta_1, \beta_2) \to (\alpha)$ | $(.5, .999)$ |
| LEARNING RATE $(f_\theta, \pi_\psi, Q_\phi)$ | $2 \times 10^{-4}$ (CHEETAH) |
| | $1 \times 10^{-3}$ (REST) |
| LEARNING RATE $(\alpha)$ | $1 \times 10^{-4}$ |
| BATCH SIZE | 512 (CHEETAH, PENDULUM) |
| | 256 (REST) |
| $Q$ FUNCTION EMA $\tau$ | 0.01 |
| CRITIC TARGET UPDATE FREQ | 2 |
| CONVOLUTION FILTERS | $[32, 32, 32, 32]$ |
| CONVOLUTION STRIDES | $[2, 2, 2, 1]$ |
| CONVOLUTION FILTER SIZE | 3 |
| ENCODER EMA $\tau$ | 0.05 |
| LATENT DIMENSION | 50 |
| HIDDEN UNITS (MLP) | $[1024, 1024]$ |
| DISCOUNT $\gamma$ | .99 |
| INITIAL TEMPERATURE | 0.1 |

*Table 2.* Atari, PPO Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| OBSERVATION RENDERING | $(84, 84)$, GREY |
| STACKED FRAMES | 4 |
| ACTION REPEAT | 4 |
| OPTIMIZER | ADAM |
| LEARNING RATE | $2.5 \times 10^{-4}$ |
| PARALLEL ENVIRONMENTS | 16 |
| SAMPLING INTERVAL | 128 |
| LIKELIHOOD RATIO CLIP, $\epsilon$ | 0.1 |
| PPO EPOCHS | 4 |
| PPO MINIBATCHES | 4 |
| CONVOLUTION FILTERS | $[32, 64, 64]$ |
| CONVOLUTION FILTER SIZES | $[8, 4, 3]$ |
| CONVOLUTION STRIDES | $[4, 2, 1]$ |
| HIDDEN UNITS (MLP) | $[512]$ |
| DISCOUNT $\gamma$ | .99 |
| GENERALIZED ADVANTAGE ESTIMATION $\lambda$ | 0.95 |
| LEARNING RATE ANNEALING | LINEAR |
| ENTROPY BONUS COEFFICIENT | 0.01 |
| EPISODIC LIVES | FALSE |
| REPEAT ACTION PROBABILITY | 0.25 |
| REWARD CLIPPING | $\pm 1$ |
| VALUE LOSS COEFFICIENT | 1.0 |

*Table 3.* DMLab, PPO Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| OBSERVATION RENDERING | $(72, 96)$, RGB |
| STACKED FRAMES | 1 |
| ACTION REPEAT | 4 |
| OPTIMIZER | ADAM |
| LEARNING RATE | $2.5 \times 10^{-4}$ |
| PARALLEL ENVIRONMENTS | 16 |
| SAMPLING INTERVAL | 128 |
| LIKELIHOOD RATIO CLIP, $\epsilon$ | 0.1 |
| PPO EPOCHS | 1 |
| PPO MINIBATCHES | 2 |
| CONVOLUTION FILTERS | $[32, 64, 64, 64]$ |
| CONVOLUTION FILTER SIZES | $[8, 4, 3, 3]$ |
| CONVOLUTION STRIDES | $[4, 2, 1, 1]$ |
| HIDDEN UNITS (LSTM) | $[256]$ |
| SKIP CONNECTIONS | CONV 3, 4; LSTM |
| DISCOUNT $\gamma$ | .99 |
| GENERALIZED ADVANTAGE ESTIMATION $\lambda$ | 0.97 |
| LEARNING RATE ANNEALING | NONE |
| ENTROPY BONUS COEFFICIENT | 0.01 (EXPLORE) |
|  | 0.0003 (LASERTAG) |
| VALUE LOSS COEFFICIENT | 0.5 |

## A.4. Online ATC Settings

*Table 4.* Common ATC Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| RANDOM SHIFT PAD | $\pm 4$ |
| LEARNING RATE | $1 \times 10^{-3}$ |
| LEARNING RATE ANNEALING | COSINE |
| TARGET UPDATE INTERVAL | 1 |
| TARGET UPDATE $\tau$ | 0.01 |
| PREDICTOR HIDDEN SIZES, $h_\psi$ | [512] |
| REPLAY BUFFER SIZE | $1 \times 10^5$ |

*Table 5.* DMControl ATC Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| RANDOM SHIFT PROBABILITY | 1 |
| BATCH SIZE | AS RL (INDIVIDUAL OBSERVATIONS) |
| TEMPORAL SHIFT, $k$ | 1 |
| MIN AGENT STEPS TO UL | $1 \times 10^4$ |
| MIN AGENT STEPS TO RL | $1 \times 10^4$ |
| UL UPDATE SCHEDULE | AS RL |
| | (2X CHEETAH) |
| LATENT SIZE | 128 |

*Table 6.* Atari ATC Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| RANDOM SHIFT PROBABILITY | 0.1 |
| BATCH SIZE | 512 (32 TRAJECTORIES OF 16 TIME STEPS) |
| TEMPORAL SHIFT, $k$ | 3 |
| MIN AGENT STEPS TO UL | $5 \times 10^4$ |
| MIN AGENT STEPS TO RL | $1 \times 10^5$ |
| UL UPDATE SCHEDULE | ANNEALED QUADRATICALLY FROM 6 PER SAMPLER ITERATION |
| | ($1 \times 10^4$ ONCE AT $1 \times 10^5$ STEPS FOR WEIGHT INITIALIZATION) |
| LATENT SIZE | 256 |

*Table 7.* DMLab ATC Hyperparameters.

| HYPERPARAMETER | VALUE |
|---|---|
| RANDOM SHIFT PROBABILITY | 1 |
| BATCH SIZE | 512 (INDIVIDUAL OBSERVATIONS) |
| TEMPORAL SHIFT, $k$ | 3 |
| MIN AGENT STEPS TO UL | $5 \times 10^4$ |
| MIN AGENT STEPS TO RL | $1 \times 10^5$ |
| UL UPDATE SCHEDULE | 2 PER SAMPLER ITERATION |
| LATENT SIZE | 256 |

## A.5. Offline Pre-Training Details

We conducted coarse hyperparameter sweeps to tune each competing UL algorithm. In all cases, the best setting is the one shown in our comparisons.

When our VAEs include a time difference between input and reconstruction observations, we include one hidden layer with action additionally input between the encoder and decoder. We tried both 1.0 and 0.1 KL-divergence weight in the VAE loss, and found 0.1 to perform better in both DMControl and Atari.

**DMControl** For the VAE, we experimented with 0 and 1 time step difference between input and reconstruction target observations and training for either $1 \times 10^4$ or $5 \times 10^4$ updates. The best settings were 1-step temporal, and $5 \times 10^4$ updates, with batch size 128. ATC used 1-step temporal, $5 \times 10^4$ updates (although this can be significantly decreased), and batch size 256 (including CHEETAH). The pretraining data set consisted of the first $5 \times 10^4$ transitions from a RAD-SAC agent learning each task, including $5 \times 10^3$ random actions. Within this span, CARTPOLE and BALL_IN_CUP learned completely, but WALKER and CHEETAH reached average returns of 514 and 630, respectively (collected without the compressive convolution).

**DMLab** For Pixel Control, we used the settings from (Hessel et al., 2019) (see the appendix therein), except we used only empirical returns, computed offline (without bootstrapping). For CPC, we tried training batch shapes, $batch \times time$ in (64, 8), (32, 16), (16, 32), and found the setting with rollouts of length 16 to be best. We contrasted all elements of the batch against each other, rather than only forward constrasts. In all cases we also used 16 steps to warmup the LSTM. For all algorithms we tried learning rates $3 \times 10^{-4}$ and $1 \times 10^{-3}$ and both $5 \times 10^4$ and $1.5 \times 10^5$ updates. For ATC and CPC, the lower learning rate and higher number of updates helped in LASERTAG especially. The pretraining data was $125 \times 10^3$ samples from partially trained RL agents receiving average returns of 127 and 6 in EXPLORE_GOAL_LOCATIONS_SMALL and LASERTAG_THREE_OPPONENTS_SMALL, respectively.

**Atari** For the VAE, we experimented with 0, 1, and 3 time step difference between input and reconstruction target, and found 3 to work best. For ST-DIM we experimented with 1, 3, and 4 time steps differences, and batch sizes from 64 to 256, learning rates $1 \times 10^{-3}$ and $5 \times 10^{-4}$. Likewise, 3-step delay worked best. For the inverse model, we tried 1- and 3-step predictions, with 1-step working better overall, and found random shift augmentation to help. For pixel control, we used the settings in (Jaderberg et al., 2017), again with full empirical returns. We ran each algorithm for up to $1 \times 10^5$ updates, although final ATC results used $5 \times 10^4$ updates. We ran each RL agent with and without observation normalization on the latent image and observed no difference in performance. Pretraining data was $125 \times 10^3$ samples sourced from the replay buffer of DQN agents trained for $15 \times 10^6$ steps with epsilon-greedy $\epsilon = 0.1$. Evaluation scores were:

*Table 8.* Atari Pre-Training Data Source Agents.

| GAME | EVALUATION SCORE |
|---:|:---|
| ALIEN | $1,800$ |
| BREAKOUT | $279$ |
| FROSTBITE | $1,400$ |
| GRAVITAR | $390$ |
| PONG | $18$ |
| QBERT | $8,800$ |
| SEAQUEST | $11,000$ |
| SPACE INVADERS | $1,200$ |