
Beyond Fine-Tuning: Transferring Behavior in Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Designing agents that acquire knowledge autonomously and use it to solve new
2 tasks efficiently is an important challenge in reinforcement learning. Knowledge
3 acquired during an unsupervised pre-training phase is often transferred by fine-
4 tuning neural network weights once rewards are exposed, as is common practice
5 in supervised domains. Given the nature of the reinforcement learning problem,
6 we argue that standard fine-tuning strategies alone are not enough for efficient
7 transfer in challenging domains. We introduce *Behavior Transfer* (BT), a technique
8 that leverages pre-trained policies for exploration and that is complementary to
9 transferring neural network weights. Our experiments show that, when combined
10 with large-scale pre-training in the absence of rewards, existing intrinsic motivation
11 objectives can lead to the emergence of complex behaviors. These pre-trained
12 policies can then be leveraged by BT to discover better solutions than without
13 pre-training, and combining BT with standard fine-tuning strategies results in
14 additional benefits. The largest gains are generally observed in domains requiring
15 structured exploration, including settings where the behavior of the pre-trained
16 policies is misaligned with the downstream task.

17 1 Introduction

18 Transfer in deep learning is often performed through parameter initialization followed by fine-tuning,
19 a technique that allows to leverage the power of deep networks in domains where labelled data
20 is scarce [64, 16, 65, 22, 15]. This builds on the intuition that the pre-trained model will map
21 inputs to a feature space where the downstream task is easy to perform. When combined with
22 methods that can leverage massive amounts of unlabelled data for pre-training, this transfer strategy
23 has led to unprecedented results in domains like computer vision [31, 30] and natural language
24 processing [15, 54]. The success of these approaches has led to an ever-growing interest in developing
25 techniques for pre-training large scale models on unlabelled data [9, 13, 24].

26 In the reinforcement learning (RL) context, unsupervised methods that learn in the absence of reward
27 have also garnered much research attention [23, 21, 49, 19, 29]. The benefits of unsupervised pre-
28 training are typically evaluated by their ability to enable efficient transfer to previously unseen reward
29 functions [28]. In spite of their different approaches to unsupervised RL, most of the top-performing
30 methods in this setting transfer knowledge through neural network weights. Such approaches deal
31 with the data inefficiency associated to training neural networks with gradient descent, similarly to
32 what is done in supervised learning, e.g. by pre-training encoders that extract representations from
33 observations [63]. However, RL introduces a challenge that is not present in supervised learning: the
34 agent is responsible for collecting the right data to learn from. This introduces a second source of
35 inefficiency from which transfer approaches can also suffer if they rely on unstructured exploration
36 strategies after pre-training, as these can lead to exponentially larger data requirements in complex

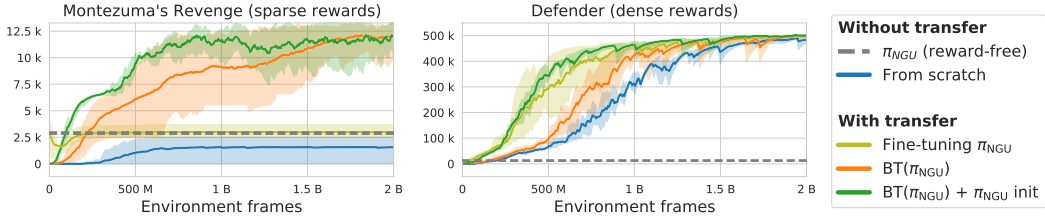


Figure 1: Comparison of transfer strategies on Montezuma’s Revenge and Defender after pre-training a policy with NGU [52] in the absence of reward. The benefits of our proposed approach to leverage pre-trained behavior for exploration, Behavior Transfer (BT), are complementary to the gains provided by pre-trained weight initialization followed by fine-tuning.

37 downstream environments [48, 47]. To address this problem, one could consider fine-tuning policies
 38 that produce meaningful behavior [46, 56], but this approach quickly disregards the pre-trained
 39 behavior when learning in the downstream task due to catastrophic forgetting.

40 In this work, we explicitly separate the transfer of behaviour and weights. We propose to make
 41 use of the pre-trained behaviour itself (i.e., the pre-trained policy mapping from observations to
 42 actions) in contrast to pre-trained neural network weights for further fine-tuning. While pre-trained
 43 behavior has been used before for *exploitation* [5, 60, 2, 3], our approach employs pre-trained policies
 44 to aid with *exploration* as well to collect experience that can be leveraged via off-policy learning.
 45 This strategy accelerates learning, as the agent is exposed to potentially useful experience earlier in
 46 training, without compromising the quality of the discovered solution when the pre-trained behavior
 47 is not aligned with the downstream task. We expose the pre-trained behaviour to the downstream
 48 agent in two ways: firstly, as an extra exploratory strategy that, when randomly activated, persists for
 49 a number of steps, and secondly as an additional pseudo-action for the learned value function where
 50 the agent may elect to defer action selection to the pre-trained policy instead of choosing itself. We
 51 call this approach Behavior Transfer (BT).

52 Defining unsupervised RL objectives remains an open problem, and solutions are generally influenced
 53 by how the acquired knowledge will be used for solving downstream tasks. Instead of proposing yet
 54 another objective for unsupervised pre-training, we turn to existing techniques for training policies in
 55 the absence of reward and make our choice based on two general requirements. First, the objective
 56 should scale gracefully with increased compute and data. This has been key for the success of
 57 self-supervised approaches in other domains [9, 36], and we argue that it is an important property for
 58 unsupervised RL as well. Second, the pre-training stage should return a policy that produces complex
 59 behavior that may be leveraged in a subsequent transfer stage. The *Never Give Up* (NGU) [52]
 60 intrinsic reward meets both requirements, and our experiments show that large-scale pre-training with
 61 this objective leads to state of the art scores in the reward-free Atari benchmark.

62 Figure 1 exemplifies our main findings. We pre-train behaviour using the intrinsic NGU reward during
 63 a long unsupervised phase without rewards. This gives rise to exploratory behaviors that seek to visit
 64 many different states throughout an episode, and we then compare different strategies for leveraging
 65 the acquired knowledge once rewards are reinstated. While fine-tuning the pre-trained weights
 66 enables faster learning, the exploratory behavior of the pre-trained policy is quickly disregarded as it
 67 is exposed to rewards. On the other hand, Behavior Transfer (BT) does not modify the pre-trained
 68 policy while learning in the new task and is able to achieve higher end scores thanks to better
 69 exploration. These two strategies are not mutually exclusive, and BT also benefits from the faster
 70 convergence provided by initializing neural networks with pre-trained weights when these encode
 71 useful information for solving the downstream task.

72 Our contributions can be summarized as follows. (1) We propose *Behavior Transfer* (BT), a technique
 73 that leverages pre-trained policies for exploration by treating them as black boxes that are not modified
 74 during learning on the downstream task. BT uses the pre-trained policy to collect experience in
 75 two ways, namely randomly-triggered temporally-extended exploration and one-step calls based on
 76 value estimates. (2) Our experiments show that large-scale unsupervised pre-training with existing
 77 intrinsic rewards can produce meaningful behavior, achieving state of the art results in the reward-free
 78 Atari benchmark. These results suggest that scale is key for unsupervised RL, akin to what has been
 79 observed in supervised settings. (3) We provide extensive empirical evidence demonstrating the

80 benefits of leveraging pre-trained behavior via BT. Our approach obtains the largest gains in hard
 81 exploration games, where it almost doubles the median human normalized score achieved by our
 82 strongest baseline. Furthermore, we show that BT is able to leverage a single task-agnostic policy
 83 to solve multiple tasks in the same environment and to achieve high performance even when the
 84 pre-trained policies are misaligned with the task being solved. (4) BT brings benefits to the table
 85 that are complementary to those provided by reusing pre-trained neural network weights, and we
 86 empirically show that combining these two strategies can result in larger gains.

87 2 Preliminaries

88 The interaction between the agent and the environment is modelled as a Markov Decision Pro-
 89 cess (MDP) [53]. An MDP is defined by the tuple $(\mathcal{S}, \mathcal{A}, P, d_0, R, \gamma)$ where \mathcal{S} and \mathcal{A} are the state
 90 and action spaces, $P(s'|s, a)$ is the probability of transitioning from state s to s' after taking action a ,
 91 $d_0(s)$ is the probability distribution over initial states, $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward function,
 92 and $\gamma \in [0, 1)$ is the discount factor. The goal is to find a policy $\pi(a|s)$ that maximizes the expected
 93 return, $G_t = \sum_{t=0}^{\infty} \gamma^t R_t$, where $R_t = r(S_t, A_t, S_{t+1})$. A principled way to address this problem
 94 is to use methods that compute action-value functions, $Q^\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$, where
 95 $\mathbb{E}_\pi[\cdot]$ denotes expectation over transitions induced by π [53].

96 We consider a setting where the agent is allowed to first learn within an MDP without rewards,
 97 $\mathcal{M}^R = (\mathcal{S}, \mathcal{A}, P, d_0)$, for a long period of time. The knowledge acquired during the reward-free
 98 stage is later leveraged when maximizing reward in new MDPs that share the same underlying
 99 dynamics but have different reward functions, $\mathcal{M}_i = (\mathcal{S}, \mathcal{A}, P, d_0, R_i, \gamma_i)$. Interactions between the
 100 agent and the environment are often assumed to incur a cost, but we will consider this cost to be
 101 relevant only for transitions with reward [28]. Even if the cost of unsupervised pre-training becomes
 102 non-negligible, it can be amortized when the acquired task-agnostic knowledge is leveraged to solve
 103 multiple tasks efficiently [15, 9]. Indeed, we would expect this transfer setting to become more
 104 relevant as the community moves towards more complex environments, where one may want to
 105 train agents to maximize multiple reward functions under constant dynamics. In the limit, one could
 106 consider the real world: it has constant or slowly changing dynamics, and humans are able to leverage
 107 previously acquired skills to quickly master new tasks.

108 3 Behavior Transfer

109 Transfer in supervised domains often exploits the fact that related tasks might be solved using similar
 110 representations. This practice deals with the data inefficiency of training large neural networks
 111 with stochastic gradient descent. However, there is an additional source of data inefficiency when
 112 training RL agents: unstructured exploration. Fine-tuning a pre-trained exploratory policy arises as
 113 a potential strategy for overcoming this problem, as the agent will observe rich experience much
 114 earlier in training than when initializing the policy randomly, but this approach suffers from important
 115 limitations. Learning in the downstream task can lead to catastrophically forgetting the pre-trained
 116 policy, thus prematurely disregarding its exploratory behavior. Moreover, the same neural network
 117 architecture needs to be used for both the pre-trained and the downstream policies, which in practice
 118 also imposes a limitation on the type of RL methods that can be employed in the adaptation stage (for
 119 instance, if the pre-trained policy was trained using a policy-based method, it might not be possible
 120 to fine-tune it using a value-based approach).

121 Let us assume that we have access to a pre-trained policy that exhibits exploratory behavior, and
 122 defer the discussion on how to train this policy to Section 4. Following such a policy might bring
 123 the agent to states that are unlikely to be visited with unstructured exploration techniques such as
 124 ϵ -greedy [59]. This property has the potential of accelerating learning even when the behavior of
 125 the pre-trained policy is not aligned with the downstream task, as it will effectively shorten the
 126 path between otherwise distant states [42]. Leveraging pre-trained policies for exploration differs
 127 from other approaches in the literature that use such policies directly for exploitation, e.g. via
 128 zero-shot transfer [19], methods that define a higher-level policy that alternates between the given
 129 policies [5, 60], or within the framework of generalized policy updates [4]. Exploring with pre-trained
 130 policies can accelerate convergence by providing useful experience to the agent, which is possible
 131 even when the pre-training and downstream tasks are misaligned. However, strategies that directly
 132 use the pre-trained policies for exploitation may result in sub-optimal solutions in such scenario [2].

133 We propose to leverage the behavior of pre-trained policies during transfer to aid with exploration. An
 134 explicit distinction between behavior and representation is made by considering pre-trained policies as
 135 black boxes that take observations and return actions. This strategy is agnostic to how the pre-trained
 136 behavior is encoded and is not restricted to learned policies. We rely on off-policy learning methods
 137 during transfer to leverage the behavior of a pre-trained policy $\pi_p(a|s)$. We keep π_p fixed during
 138 transfer, which prevents catastrophic forgetting of the original behavior when it is parameterized by a
 139 neural network (i.e., we instantiate and train a new policy with its own set of parameters). We propose
 140 *Behavior Transfer* (BT), which leverages two complementary strategies to achieve this. Since BT
 141 is agnostic to the method used to pre-train policies, $BT(\pi_p)$ refers to behavior being transferred
 142 from policy π_p . We formalize BT in the context of value-based Q-learning agents, although similar
 143 derivations are in principle possible for alternative off-policy learning methods. Pseudo-code for BT
 144 is provided in Algorithm 1.

145 **Temporally-extended exploration.** We draw inspiration from Lévy flights [61], a class of ecological
 146 models for animal foraging, where a fixed direction is followed for a duration sampled from a
 147 heavy-tailed distribution. This principle was implemented in the context of exploration in RL by
 148 ϵz -greedy [14], which encodes the notion of direction in the environment via exploration options that
 149 repeat the same action throughout the entire flight. Since π_p is more likely to encode a meaningful
 150 notion of direction in complex environments than action repeats, we propose a variant of ϵz -greedy
 151 where π_p is used as the exploration option. An exploratory flight might be started at any step with
 152 some probability. The duration for the flight is sampled from a heavy-tailed distribution (Zeta with
 153 $\mu = 2$ in all our experiments), and control is handed over to π_p during the complete flight. When not
 154 in a flight, actions are sampled from the behavior policy obtained while maximizing the task reward
 155 (e.g. an ϵ -greedy derived from the estimated Q values).

156 **Extra action.** The previous approach switches to π_p during experience collection blindly, and we
 157 now consider an alternative strategy for triggering these switches based on value. This can be easily
 158 implemented through an extra action which samples an action from π_p , which also allows the agent to
 159 use the pre-trained policy at test time if deemed beneficial. More formally, this amounts to training a
 160 policy over an expanded action set $\mathcal{A}^+ = \mathcal{A} \cup \{a_+\}$, where a_+ is resolved by sampling an action from
 161 π_p , $a' \sim \pi_p(s)$ (with $a' \in \mathcal{A}$). The additional action can be seen as an option that can be initiated
 162 from any state and always terminates after a single step. Note that selecting the option will lead to
 163 the same outcome as if the agent had selected a' as a primitive action, and we take advantage of this
 164 observation by using the return of following the option as target to fit both $Q(s, \pi_p(s))$ and $Q(s, a')$.
 165 Intuitively, this approach induces a bias that favours actions selected by π_p , accelerating the collection
 166 of rewarding transitions when the pre-trained policy is somewhat aligned with the downstream task.
 167 Otherwise, the agent can learn to ignore π_p as training progresses by selecting other actions.

Algorithm 1: Experience collection pseudo-code for BT

Input: Action set, \mathcal{A} ; additional action, a_+ ; extended action set, $\mathcal{A}^+ = \mathcal{A} \cup \{a_+\}$; pre-trained policy, π_p ; Q-value estimate for the current policy, $Q^\pi(s, a) \forall a \in \mathcal{A}^+$; probability of taking an exploratory action, ϵ ; probability of starting a flight, ϵ_{levy} ; flight length distribution, $\mathcal{D}(\mathbb{N})$

```

while True do
  n ← 0 // flight length
  while episode not ended do
    Observe state s
    if n == 0 and random() ≤ εlevy then n ~ D(N) // sample flight length
    if n > 0 then
      n ← n - 1
      a ~ πp(s)
    else
      if random() ≤ ε then a ~ Uniform(A+) else a ← arg maxa' ∈ A+ [Qπ(s, a')]
      if a == a+ then a ~ πp(s)
    end
    Take action a
  end
end

```

168 4 Reward-free pre-training

169 It is a common practice to derive objectives for proxy tasks in order to drive learning in the absence
170 of reward functions, and there exists a plethora of different approaches in the literature. Model-based
171 approaches can learn world models from unsupervised interaction [26]. However, the diversity of
172 the training data will impact the accuracy of the model [57] and deploying this type of approach
173 in visually complex domains like Atari remains an open problem [27]. Unsupervised RL has also
174 been explored through the lens of *empowerment* [55, 44], which studies agents that aim to discover
175 intrinsic options [23, 19]. While these options can be leveraged by hierarchical agents [21] or
176 integrated within the universal successor features framework [2, 3, 8, 28], their potential lack of
177 coverage generally limits their applicability to complex downstream tasks [12]. An alternative
178 objective is that of exploring the environment by finding policies that induce maximally entropic state
179 distributions [29, 40], although this might become extremely inefficient in high-dimensional state
180 spaces without proper priors [41, 63].

181 Recall that our goal is to devise a pre-training objective that can help reduce the amount of interaction
182 needed by the agent to collect relevant experience when learning in a downstream task. We argue that
183 such objective needs to meet two requirements. First, as suggested by results in other domains [9, 36],
184 it should scale gracefully as the amount of compute and experience used for pre-training are increased.
185 This contrasts with the training regimes used in most unsupervised RL approaches, which use a
186 relatively small amount of experience [28, 41, 63] when compared to distributed agents that do make
187 use of rewards [33, 18, 37]. Second, it must encourage the emergence of complex behaviors such as
188 navigation or manipulation skills. It has been argued that exploring the environment efficiently will
189 serve as a proxy for developing such behaviors [38], and exploration bonuses have been shown to
190 produce meaningful behavior in the absence of reward [49, 10]. However, many exploration bonuses
191 vanish over the course of training and thus may not be well-suited for a long unsupervised pre-training
192 phase. It can be shown that many intrinsic rewards aim at maximizing the entropy of all states visited
193 during training, and so the final policy does not necessarily exhibit exploratory behavior [40].

194 We propose to use Never Give Up (NGU) [52] as a means for training exploratory policies in an
195 unsupervised setting. The NGU intrinsic reward proposes a curiosity-driven approach for training
196 persistent exploratory policies which combines per-episode and life-long novelty. The per-episode
197 novelty, r_t^{episodic} , rapidly vanishes over the course of an episode, and it is designed to encourage self-
198 avoiding trajectories. It is computed by comparing a representation of the current observation, $f(s_t)$,
199 to those of all the observations visited in the current episode, $M = \{f(s_0), f(s_1), \dots, f(s_{t-1})\}$,
200 where $f : \mathcal{S} \rightarrow \mathbb{R}^p$ is an embedding function trained using a self-supervised inverse dynamics
201 model [49]. Such a mapping concentrates on the controllable aspects of the environment, ignoring
202 all the variability present in the observation that is not affected by the action taken by the agent.
203 The life-long novelty, α_t , slowly vanishes throughout training, and it is computed by using Random
204 Network Distillation (RND) [11]. With this, the intrinsic reward r_t^{NGU} is defined as follows:

$$r_t^{\text{NGU}} = r_t^{\text{episodic}} \cdot \min \{ \max \{ \alpha_t, 1 \}, L \}, \text{ with } r_t^{\text{episodic}} = \frac{1}{\sqrt{\sum_{f(s_i) \in N_k} K(f(s_t), f(s_i))} + c}} \quad (1)$$

205 where L is a fixed maximum reward scaling, N_k is the set containing the k -nearest neighbors of $f(s_t)$
206 in M , c is a constant and $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^+$ is a kernel function satisfying $K(x, x) = 1$ (which
207 can be thought of as approximating pseudo-counts [52]). The episodic component of the reward
208 in Equation 1 is reset by emptying M with each episode, thus the NGU reward does not vanish
209 throughout the training process. This makes it suitable for driving learning in task-agnostic settings.
210 Further details on NGU are reported in the supplementary material.

211 5 Experiments

212 Agents are evaluated in the Atari suite [7], a benchmark that presents a variety of challenges and that
213 is a common test ground for RL agents with unsupervised pre-training [28, 41, 56]. Experiments are
214 run using the distributed R2D2 agent [37] with 256 CPU actors and a single GPU learner. Policies
215 use the same Q-Network architecture as Agent57 [51], which is composed by a convolutional torso
216 followed by an LSTM [32] and a dueling head [62]. Hyperparameters and a detailed description of
217 the full distributed setting are provided in the supplementary material. All reported results are the
218 average over three random seeds.

219 **Reward-free learning.** The amount of task reward collected by unsupervised policies is often
 220 used as a proxy to measure their quality [19]. While the actual utility of these policies will not
 221 be revealed until they are leveraged for transfer, this proxy lets us evaluate whether the discovered
 222 behavior changes as longer pre-training budgets are allowed. We compare unsupervised NGU policies
 223 against VISR [28] and APT [41], which utilize a small amount of supervised interaction to adapt
 224 the pre-trained policies. We also consider two additional unsupervised baselines: (i) a constant
 225 positive reward at each timestep that favours long episodes, which correlate with high scores in
 226 some games [10], and (ii) RND [11], which rewards life-long novelty. Note that the RND reward
 227 vanishes, but we include it in our analysis because it was previously used by Burda et al. [10] in
 228 this setting and implementation choices such as reward normalization may prevent it from fading
 229 in practice. Figure 2 (left) shows how the zero-shot transfer performance of unsupervised policies
 230 evolves during a long pre-training phase. NGU reaches the highest scores, but both NGU and RND
 231 eventually outperform VISR and APT even though these used supervised interaction. In Table 2 of
 232 Appendix C we show that unsupervised NGU policies largely outperform several other baselines
 233 using the standard pre-training and adaptation setting. These results highlight the importance of
 234 large-scale unsupervised pre-training in RL, similarly to the trend observed in supervised domains [9].

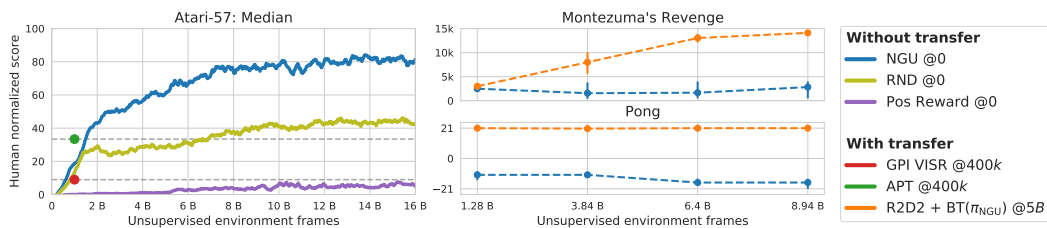


Figure 2: Performance as a function of the pre-training budget. @ N represents the number of frames with reward utilized for transfer. **(Left)** Median human normalized score across the 57 games in the Atari suite. We observe the emergence of useful behavior when optimizing an intrinsic reward during a long unsupervised pre-training of 16B frames, which contrasts with the shorter pre-training of 1B frames in previous works [28, 41]. **(Right)** Scores in the games of Montezuma’s Revenge (sparse rewards) and Pong (dense reward), before and after transfer, as a function of the pre-training budget. A longer pre-training benefits transfer in hard exploration games even if the zero-shot transfer score of the unsupervised policies does not increase.

235 **Transfer setting.** Transfer approaches are typically evaluated in the Atari benchmark with a budget
 236 of 100k RL interactions with reward (400k frames), but we propose to allow a longer adaptation
 237 phase. Randomly initialized networks tend to overfit in these very low data regimes without strong
 238 regularization [39], and we are interested in studying the impact of leveraging behavior both in
 239 isolation and combined with transfer via pre-trained weights. Moreover, since the pre-trained policies
 240 are already competent in the downstream tasks, 100k interactions are exhausted after few episodes
 241 and may be insufficient for improving performance. For these reasons, we provide results with up
 242 to 1.25B RL steps of supervised interaction (5B frames). This allows evaluating both convergence
 243 speed and asymptotic performance, while still being a relatively small budget for these distributed
 244 agents with hundreds of actors [51].

245 **Transfer via behavior.** We start by studying the impact of leveraging behavior in isolation, i.e. with-
 246 out transferring pre-trained weights, when learning in downstream tasks. We compare BT against two
 247 baselines that do not use pre-trained behavior, namely the standard R2D2 agent [37] that uses ϵ -greedy
 248 policies for exploration [59], as well as a variant of R2D2 with ϵ_Z -greedy exploration [14]. Figure 3
 249 shows that BT is superior to both baselines for any amount of environment interaction with rewards,
 250 converging faster early in training and also obtaining higher asymptotic performance. These results
 251 also demonstrate the generality of the proposed approach, as it is able to benefit from both RND
 252 and NGU policies. Note that BT performs particularly well in the set of six hard exploration games¹
 253 defined by Bellemare et al. [6], which is aligned with our intuition that reusing behavior helps over-
 254 coming the inefficiency associated to unstructured exploration. Figure 2 (right) confirms that a long
 255 pre-training phase is especially important in hard exploration games such as Montezuma’s Revenge,
 256 even if they do not translate into higher zero-shot transfer scores, as it produces more exploratory be-
 257 havior. On the other hand, the performance after transfer is independent of the amount of pre-training
 258 in dense reward games like Pong, where unstructured exploration is enough to reach optimal scores.

¹gravitar, montezuma_revenge, pitfall, private_eye, solaris, venture

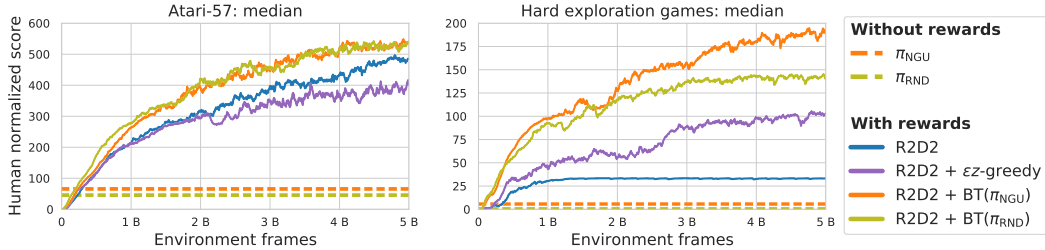


Figure 3: Median human normalized scores for R2D2-based agents trained from scratch. **(Left)** Full Atari suite. **(Right)** Subset of hard exploration games.

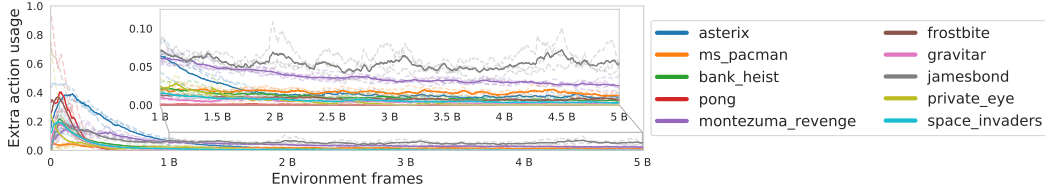


Figure 4: Usage of the extra action in $\text{BT}(\pi_{\text{NGU}})$, computed as the fraction of steps within an episode in which it is selected by the agent. The usage peaks early in training and slowly decreases afterwards as the new policy becomes stronger at the task.

259 **Ablation studies.** In order to gain insight on each of the components in BT, we run experiments
 260 on a subset of 12 games² requiring different amounts of exploration and featuring both dense and
 261 sparse rewards. $\text{BT}(\pi_{\text{NGU}})$ achieves a median score of 368 in this subset, which compares favorably
 262 to the 196 median score of R2D2 with ϵ -greedy exploration. Removing either the extra action or
 263 the temporally-extended exploration reduces the median score of $\text{BT}(\pi_{\text{NGU}})$ to 224. These results
 264 suggest that the gains provided by both strategies are complementary, and both are responsible for the
 265 strong performance of BT. To provide further insight about the benefits of BT, Figure 4 reports the
 266 fraction of steps per episode in which the extra action is selected by the greedy policy. It hints at the
 267 emergence of a schedule over the usage of the pre-trained policy, which increases early in training
 268 and decays afterwards. We hypothesize that this is due to the fact that the unsupervised policies
 269 obtain large episodic returns, but their behavior is suboptimal when maximizing discounted rewards.
 270 These policies take many exploratory actions in between rewards, and so the agent eventually figures
 271 out more efficient strategies for reaching rewarding states by using primitive actions.

272 **Transfer to multiple tasks.** An appealing property of task-agnostic knowledge is that it can be
 273 leveraged to solve multiple tasks. In the RL setting, this can be evaluated by leveraging a single
 274 task-agnostic policy for solving multiple tasks (i.e. reward functions) in the same environment. We
 275 evaluate whether the unsupervised NGU policies can be useful beyond the standard Atari tasks by
 276 creating two alternative versions of Ms Pacman and Hero with different levels of difficulty. The
 277 goal in the modified version of Ms Pacman is to eat vulnerable ghosts, with pac-dots giving 0 (easy
 278 version) or -10 (hard version) points. In the modified version of Hero, saving miners gives a fixed
 279 return of 1000 points and dynamiting walls gives either 0 (easy version) or -300 (hard version) points.
 280 The rest of rewards are removed, e.g. eating fruit in Ms Pacman or the bonus for unused power units
 281 in Hero. Note that even in the easy version of the games exploration is harder than in their original
 282 counterparts, as there are no small rewards guiding the agent towards its goals. Exploration is even
 283 more challenging in the hard version of the games, as the intermediate rewards work as a deceptive
 284 signal that takes the agent away from its actual goal. In this case, finding rewarding behaviors requires
 285 a stronger commitment to an exploration strategy. Unsupervised NGU policies often achieve very low
 286 or even negative rewards in this setting, which contrasts with the strong performance they showed
 287 when evaluated under the standard game reward. Figure 5 shows that leveraging the behavior of
 288 pre-trained exploration policies provides important gains even in this adversarial scenario. These
 289 results suggest that the strong performance observed under the standard game rewards is not due to an

²Obtained by combining games used to tune hyperparameters in [28] with games where ϵ -greedy provides clear gains over ϵ -greedy as per [14]: asterix, bank_heist, frostbite, gravitar, jamesbond, montezuma_revenge, ms_pacman, pong, private_eye, space_invaders, tennis, up_n_down.



Figure 5: Scores in Atari games with modified reward functions. We train a single task-agnostic policy per environment, and leverage it to solve three different tasks: the standard game reward, a task with sparse rewards (easy), and a variant of the same task with deceptive rewards (hard).

290 alignment between the NGU reward and the game goals, but due to an efficient usage of pre-trained
 291 exploration policies.

292 **Combining pre-trained behavior and weights.** Our last batch of experiments focuses on studying
 293 transfer via pre-trained weights and its compatibility with BT. Policies are composed of a convo-
 294 lutional torso, an LSTM, and a dueling head. We consider two initialization strategies: a *partial*
 295 *initialization* approach that loads the torso and the LSTM, but initializes the head randomly; and a
 296 *full initialization* scheme where all weights are loaded. The former can be understood as transferring
 297 learned representations [63], but deferring exploration to a random policy. On the other hand, the
 298 full initialization approach can be seen as directly transferring the policy and is usually referred to as
 299 fine-tuning the pre-trained policy [46, 41, 56]. Note that these approaches only change how weights
 300 are initialized *before* training. As in previous experiments, all parameters in the new policy are trained
 301 and π_p is kept fixed when using BT. Figure 6 (top) compares agents with and without BT for different
 302 amounts of transfer via weights on the Atari benchmark. Loading pre-trained weights results in faster
 303 learning early in training, both with and without BT. The largest gains are observed in dense reward
 304 games, which translates into higher median scores across the full suite because most games belong
 305 to this category. Weights alone are not enough in hard exploration games, where leveraging the
 306 pre-trained policy via BT provides clear benefits. Perhaps surprisingly, we observe that transferring
 307 representations outperforms fine-tuning the pre-trained policy, and we hypothesize that the former
 308 is more robust to misalignments between the pre-trained policy and the downstream task. This
 309 intuition is further supported by the experiments on games with modified reward functions reported
 310 in Figure 6 (middle & bottom), where the faster learning provided by pre-trained weights often comes
 311 at the cost of lower end scores. On the other hand, BT is crucial in tasks with sparse and deceptive
 312 rewards and also benefits from pre-trained weights in tasks where positive transfer is observed.

313 6 Related work

314 Our work uses the experimental methodology presented by Hansen et al. [28]. Whereas that work only
 315 considered a fast, simplified adaptation process that limited the final performance on the downstream
 316 task, we focus on the more general case of using a previously trained policy to aid in solving the
 317 full RL problem. Hansen et al. [28] use successor features to identify which of the pre-trained tasks
 318 best matches the true reward structure, which has previously been shown to work well for multi-task
 319 transfer [3]. Bagot et al. [1] augments an agent with the ability to utilize another policy, which is
 320 learned in tandem based on an intrinsic reward function. This promising direction is complementary
 321 to our work, as it handles the case wherein there is no unsupervised pre-training phase.

322 Gupta et al. [25] provides an alternative method to meta-learn a solver for reinforcement learning prob-
 323 lems from unsupervised reward functions. This method utilizes gradient-based meta-learning [20],
 324 which makes the adaptation process standard reinforcement learning updates. This means that even if
 325 the downstream reward is far outside of the training distribution, final performance would not neces-
 326 sarily be affected. However, these methods are hard to scale to the larger networks considered here,
 327 and followup work [34] changed to memory-based meta-learning [17] which relies on information
 328 about rewards staying in the recurrent state. This makes it unsuitable to the sort of hard exploration
 329 problem our method excels at. Recent work has shown success in transferring representations learned
 330 in an unsupervised setting to reinforcement learning tasks [58]. Our representation transfer experi-

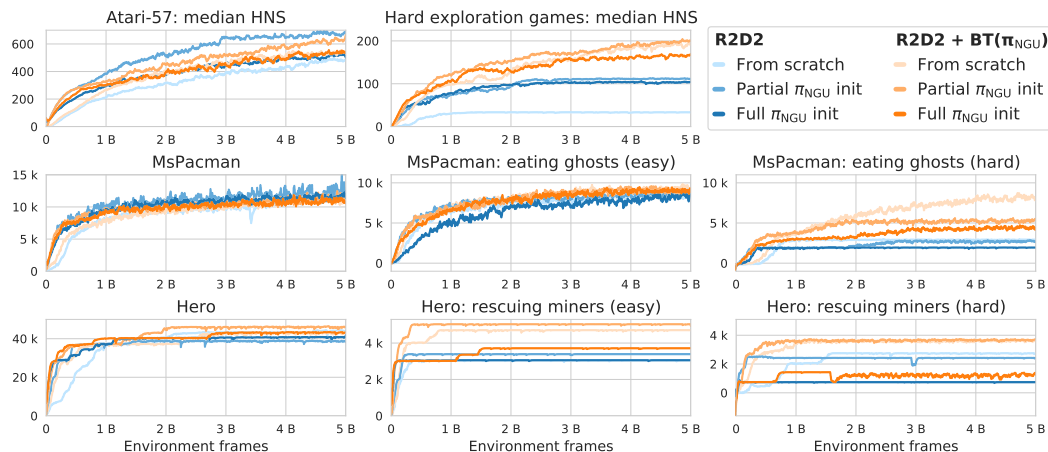


Figure 6: Performance of R2D2-based agents with different amounts of transfer via weights. Policies are composed of a CNN encoder followed by an LSTM and a dueling head. We compare training from scratch, loading all weights (Full π_{NGU} init) or all weights except those in the dueling head (Partial π_{NGU} init). **(Top)** Median human normalized scores (HNS) in the full Atari suite (left) and the subset of hard exploration games (right). **(Middle & Bottom)** Games with modified reward functions as in Figure 5.

331 ments suggest that this might handicap final performance, but the possibility also exists that different
 332 unsupervised objectives should be used for representation transfer and policy transfer.

333 7 Discussion

334 We studied the problem of transferring pre-trained behavior for exploration in reinforcement learning,
 335 an approach that is complementary to the common practice of transferring neural network weights.
 336 Our proposed approach, Behavior Transfer (BT), relies on the pre-trained policy for collecting
 337 experience in two different ways: (i) through temporally-extended exploration, which can be triggered
 338 with some probability at any step, and (ii) via one-step calls to the pre-trained policy based on value
 339 estimates. BT results in strong transfer performance when combined with exploratory policies pre-
 340 trained in the absence of reward, with the most important gains being observed in hard exploration
 341 tasks. These benefits are not due to an alignment between our pre-training and downstream tasks,
 342 as we also observed positive transfer in games where the pre-trained policy obtained low scores.
 343 In order to provide further evidence for this claim, we designed alternative tasks for Atari games
 344 involving hard exploration and deceptive rewards. Our transfer strategy outperformed all considered
 345 baselines in these settings, even when the pre-trained policy obtained very low or even negative scores,
 346 demonstrating the generality of the method. Besides disambiguating the role of the alignment between
 347 pre-training and downstream tasks, these experiments demonstrate the utility of a single task-agnostic
 348 policy for solving multiple tasks in the same environment. Finally, we also demonstrated that BT can
 349 be combined with transfer via neural network weights to provide further gains.

350 Our experimental results highlight the importance of scale when training RL agents in reward-free
 351 settings, which is one of the key factors behind the recent success of unsupervised approaches in other
 352 domains. This contrasts with the small budgets considered for reward-free RL in previous works and
 353 motivates further research in unsupervised RL approaches that scale with increased data and compute.
 354 We argue that scale is one of the missing components in reward-free RL, and it will be a necessary
 355 condition to unfold its full potential. Beyond improving the unsupervised learning phase, we are also
 356 excited about the possibilities unlocked by BT and that are not possible when transferring knowledge
 357 through weights, such as leveraging multiple pre-trained policies and deploying BT in continual
 358 learning scenarios where the agent never stops learning and keeps accumulating knowledge and skills.
 359 Future work should also study improved mechanisms for handing over control to pre-trained policies,
 360 as well as prioritizing the usage of certain behaviors over others when multiple such policies are
 361 available to the agent. This could overcome one of the current limitations of BT, which assumes that
 362 flights can be started from any state and still produce meaningful behavior.

363 References

- 364 [1] Louis Bagot, Kevin Mets, and Steven Latré. Learning intrinsically motivated options to stimulate
365 policy exploration. In *ICML Workshop on LifeLong Learning*, 2020.
- 366 [2] André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt,
367 and David Silver. Successor features for transfer in reinforcement learning. In *NeurIPS*, 2017.
- 368 [3] Andre Barreto, Diana Borsa, John Quan, Tom Schaul, David Silver, Matteo Hessel, Daniel
369 Mankowitz, Augustin Zidek, and Remi Munos. Transfer in deep reinforcement learning using
370 successor features and generalised policy improvement. In *ICML*, 2018.
- 371 [4] André Barreto, Shaobo Hou, Diana Borsa, David Silver, and Doina Precup. Fast reinforcement
372 learning with generalized policy updates. *Proceedings of the National Academy of Sciences*,
373 2020.
- 374 [5] Andrew G Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement
375 learning. *Discrete event dynamic systems*, 2003.
- 376 [6] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi
377 Munos. Unifying count-based exploration and intrinsic motivation. In *NeurIPS*, 2016.
- 378 [7] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning
379 environment: An evaluation platform for general agents. *Journal of Artificial Intelligence*
380 *Research*, 2013.
- 381 [8] Diana Borsa, André Barreto, John Quan, Daniel Mankowitz, Rémi Munos, Hado van Hasselt,
382 David Silver, and Tom Schaul. Universal successor features approximators. In *ICLR*, 2019.
- 383 [9] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
384 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
385 few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- 386 [10] Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A Efros.
387 Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*, 2018.
- 388 [11] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random
389 network distillation. *arXiv preprint arXiv:1810.12894*, 2018.
- 390 [12] Víctor Campos, Alexander Trott, Caiming Xiong, Richard Socher, Xavier Giro-i Nieto, and
391 Jordi Torres. Explore, discover and learn: Unsupervised discovery of state-covering skills. In
392 *ICML*, 2020.
- 393 [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework
394 for contrastive learning of visual representations. In *ICML*, 2020.
- 395 [14] Will Dabney, Georg Ostrovski, and André Barreto. Temporally-extended ϵ -greedy exploration.
396 In *ICLR*, 2021.
- 397 [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
398 deep bidirectional transformers for language understanding. In *NAACL*, 2019.
- 399 [16] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor
400 Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*,
401 2014.
- 402 [17] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL^2 :
403 Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*,
404 2016.
- 405 [18] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward,
406 Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. IMPALA: Scalable distributed
407 deep-RL with importance weighted actor-learner architectures. In *ICML*, 2018.

- 408 [19] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you
409 need: Learning skills without a reward function. In *ICLR*, 2019.
- 410 [20] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adapta-
411 tion of deep networks. *arXiv preprint arXiv:1703.03400*, 2017.
- 412 [21] Carlos Florensa, Yan Duan, and Pieter Abbeel. Stochastic neural networks for hierarchical
413 reinforcement learning. In *ICLR*, 2017.
- 414 [22] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for
415 accurate object detection and semantic segmentation. In *CVPR*, 2014.
- 416 [23] Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv*
417 *preprint arXiv:1611.07507*, 2016.
- 418 [24] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena
419 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi
420 Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv*
421 *preprint arXiv:2006.07733*, 2020.
- 422 [25] Abhishek Gupta, Benjamin Eysenbach, Chelsea Finn, and Sergey Levine. Unsupervised
423 meta-learning for reinforcement learning. *arXiv preprint arXiv:1806.04640*, 2018.
- 424 [26] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In
425 *NeurIPS*, 2018.
- 426 [27] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control:
427 Learning behaviors by latent imagination. In *ICLR*, 2019.
- 428 [28] Steven Hansen, Will Dabney, Andre Barreto, Tom Van de Wiele, David Warde-Farley, and
429 Volodymyr Mnih. Fast task inference with variational intrinsic successor features. In *ICLR*,
430 2020.
- 431 [29] Elad Hazan, Sham M Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum
432 entropy exploration. In *ICML*, 2019.
- 433 [30] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
434 unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.
- 435 [31] Olivier J Hénaff, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient
436 image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019.
- 437 [32] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
438 1997.
- 439 [33] Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado
440 Van Hasselt, and David Silver. Distributed prioritized experience replay. *arXiv preprint*
441 *arXiv:1803.00933*, 2018.
- 442 [34] Allan Jabri, Kyle Hsu, Abhishek Gupta, Ben Eysenbach, Sergey Levine, and Chelsea Finn. Un-
443 supervised curricula for visual meta-reinforcement learning. In *Advances in Neural Information*
444 *Processing Systems*, pages 10519–10531, 2019.
- 445 [35] Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell, Konrad
446 Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, et al. Model-
447 based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.
- 448 [36] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,
449 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
450 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 451 [37] Steven Kapturowski, Georg Ostrovski, John Quan, Remi Munos, and Will Dabney. Recurrent
452 experience replay in distributed reinforcement learning. In *ICLR*, 2019.

- 453 [38] Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time.
454 *Machine learning*, 2002.
- 455 [39] Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing
456 deep reinforcement learning from pixels. In *ICLR*, 2021.
- 457 [40] Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric Xing, Sergey Levine, and Rus-
458 lan Salakhutdinov. Efficient exploration via state marginal matching. *arXiv preprint*
459 *arXiv:1906.05274*, 2019.
- 460 [41] Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. *arXiv*
461 *preprint arXiv:2103.04551*, 2021.
- 462 [42] Yao Liu and Emma Brunskill. When simple exploration is sample efficient: Identifying sufficient
463 conditions for random exploration to yield pac rl algorithms. *arXiv preprint arXiv:1805.09045*,
464 2018.
- 465 [43] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
466 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
467 Human-level control through deep reinforcement learning. *Nature*, 2015.
- 468 [44] Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for
469 intrinsically motivated reinforcement learning. In *NeurIPS*, 2015.
- 470 [45] Rémi Munos, Tom Stepleton, Anna Harutyunyan, and Marc Bellemare. Safe and efficient
471 off-policy reinforcement learning. In *NeurIPS*, 2016.
- 472 [46] Mirco Mutti, Lorenzo Pratissoli, and Marcello Restelli. Task-agnostic exploration via policy
473 gradient of a non-parametric state entropy estimate. In *AAAI*, 2021.
- 474 [47] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via
475 bootstrapped dqn. *arXiv preprint arXiv:1602.04621*, 2016.
- 476 [48] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized
477 value functions. In *ICML*, 2016.
- 478 [49] Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration
479 by self-supervised prediction. In *ICML*, 2017.
- 480 [50] Jing Peng and Ronald J Williams. Incremental multi-step q-learning. In *Machine Learning*
481 *Proceedings 1994*. Elsevier, 1994.
- 482 [51] Adrià Puigdomènech Badia, Bilal Piot, Steven Kapturowski, Pablo Sprechmann, Alex Vitvitskyi,
483 Daniel Guo, and Charles Blundell. Agent57: Outperforming the atari human benchmark. In
484 *ICML*, 2020.
- 485 [52] Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, Bilal Piot, Steven
486 Kapturowski, Olivier Tieleman, Martín Arjovsky, Alexander Pritzel, Andrew Bolt, et al. Never
487 give up: Learning directed exploration strategies. In *ICLR*, 2020.
- 488 [53] Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*.
489 John Wiley & Sons, Inc., 1994.
- 490 [54] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
491 Language models are unsupervised multitask learners. *OpenAI Blog*, 2019.
- 492 [55] Christoph Salge, Cornelius Glackin, and Daniel Polani. Empowerment – an introduction. In
493 *Guided Self-Organization: Inception*. Springer, 2014.
- 494 [56] Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin,
495 R Devon Hjelm, Philip Bachman, and Aaron Courville. Pretraining reward-free representations
496 for data-efficient reinforcement learning. In *Self-Supervision for Reinforcement Learning*
497 *Workshop - ICLR 2021*, 2021. URL <https://openreview.net/forum?id=o5z9Le5drua>.

- 498 [57] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak
499 Pathak. Planning to explore via self-supervised world models. In *ICML*, 2020.
- 500 [58] Adam Stooke, Kimin Lee, Pieter Abbeel, and Michael Laskin. Decoupling representation
501 learning from reinforcement learning. *arXiv preprint arXiv:2009.08319*, 2020.
- 502 [59] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,
503 2018.
- 504 [60] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A
505 framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 1999.
- 506 [61] Gandhimohan M Viswanathan, V Afanasyev, SV Buldyrev, EJ Murphy, PA Prince, and H Eu-
507 gene Stanley. Lévy flight search patterns of wandering albatrosses. *Nature*, 1996.
- 508 [62] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas.
509 Dueling network architectures for deep reinforcement learning. In *ICML*, 2016.
- 510 [63] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement learning with
511 prototypical representations. In *ICML*, 2021.
- 512 [64] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in
513 deep neural networks? *arXiv preprint arXiv:1411.1792*, 2014.
- 514 [65] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In
515 *ECCV*, 2014.

516 Checklist

517 The checklist follows the references. Please read the checklist guidelines carefully for information on
518 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
519 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
520 the appropriate section of your paper or providing a brief inline description. For example:

- 521 • Did you include the license to the code and datasets? **[No]** The code and the data are
522 proprietary.

523 Please do not modify the questions and only use the provided macros for your answers. Note that the
524 Checklist section does not count towards the page limit. In your paper, please delete this instructions
525 block and only keep the Checklist section heading above along with the questions/answers below.

526 1. For all authors...

- 527 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
528 contributions and scope? **[Yes]**
- 529 (b) Did you describe the limitations of your work? **[Yes]**
- 530 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 531 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
532 them? **[Yes]**

533 2. If you are including theoretical results...

- 534 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 535 (b) Did you include complete proofs of all theoretical results? **[N/A]**

536 3. If you ran experiments...

- 537 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
538 imental results (either in the supplemental material or as a URL)? **[No]** We did not
539 include source code because it relies on non-public libraries that are specific to our
540 distributed hardware setting. However, we include all the details needed to replicate
541 our experiments.

- 542 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
543 were chosen)? [Yes]
- 544 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
545 ments multiple times)? [Yes] All our experiments were run with three different random
546 seeds. Plots report mean, min and max results. Tables report mean and standard
547 deviation.
- 548 (d) Did you include the total amount of compute and the type of resources used (e.g., type
549 of GPUs, internal cluster, or cloud provider)? [Yes]
- 550 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 551 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 552 (b) Did you mention the license of the assets? [N/A]
- 553 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 554
- 555 (d) Did you discuss whether and how consent was obtained from people whose data you're
556 using/curating? [N/A]
- 557 (e) Did you discuss whether the data you are using/curating contains personally identifiable
558 information or offensive content? [N/A]
- 559 5. If you used crowdsourcing or conducted research with human subjects...
- 560 (a) Did you include the full text of instructions given to participants and screenshots, if
561 applicable? [N/A]
- 562 (b) Did you describe any potential participant risks, with links to Institutional Review
563 Board (IRB) approvals, if applicable? [N/A]
- 564 (c) Did you include the estimated hourly wage paid to participants and the total amount
565 spent on participant compensation? [N/A]

566 **A Pseudo-code**

567 Algorithm 2 provides pseudo-code for the flight logic that controls how the pre-trained policy is
568 used for temporally-extended exploration. At each step, a flight is started with probability ϵ_{levy} .
569 The duration of the flight is sampled from a heavy-tailed distribution, $\mathcal{D}(\mathbb{N})$, similarly to ϵz -greedy
570 (c.f. Appendix B for more details). When not in a flight, the exploitative policy that maximizes the
571 extrinsic reward is derived from the estimated Q-values using the ϵ -greedy operator. This ensures that
572 all state-action pairs will be visited given enough time, as exploring only with π_p does not guarantee
573 such property.

574 Algorithm 3 provides pseudo-code for the actor logic when using the augmented action set, $\mathcal{A}^+ =$
575 $\mathcal{A} \cup \{\pi_p(s)\}$. It derives an ϵ -greedy policy over $|\mathcal{A}| + 1$ actions, where the $(|\mathcal{A}| + 1)$ -th action is
576 resolved by sampling from $\pi_p(s)$.

Algorithm 2: Experience collection pseudo-code for BT with temporally-extended exploration

Input: Action set \mathcal{A}
Input: Q-value estimate for the current policy, $Q^\pi(s, a) \forall a \in \mathcal{A}$
Input: Pre-trained policy, π_p
Input: Probability of starting a flight, ϵ_{levy}
Input: Flight length distribution, $\mathcal{D}(\mathbb{N})$

```
while True do
   $n \leftarrow 0$  // flight length
  while episode not ended do
    Observe state  $s$ 
    if  $n == 0$  and  $\text{random}() \leq \epsilon_{\text{levy}}$  then
       $n \sim \mathcal{D}(\mathbb{N})$  // sample from distribution over lengths
    end
    if  $n > 0$  then
       $n \leftarrow n - 1$ 
       $a \sim \pi_p(s)$ 
    else
       $a \sim \epsilon\text{-greedy}[Q^\pi(s, a)]$ 
    end
    Take action  $a$ 
  end
end
```

Algorithm 3: Experience collection pseudo-code for BT with an extra action

Input: Action set \mathcal{A}
Input: Additional action, a_+
Input: Extended action set, $\mathcal{A}^+ = \mathcal{A} \cup \{a_+\}$
Input: Pre-trained policy, π_p
Input: Q-value estimate for the current policy, $Q^\pi(s, a) \forall a \in \mathcal{A}^+$
Input: Probability of taking an exploratory action, ϵ

```
while True do
  while episode not ended do
    Observe state  $s$ 
    if  $\text{random}() \leq \epsilon$  then
       $a \sim \text{Uniform}(\mathcal{A}^+)$ 
    else
       $a \leftarrow \arg \max_{a' \in \mathcal{A}^+} [Q^\pi(s, a')]$ 
    end
    if  $a == a_+$  then
       $a \sim \pi_p(s)$ 
    end
    Take action  $a$ 
  end
end
```

577 **B Hyperparameters**

578 All policies use the same Q-Network architecture as Agent57 [51], which is composed by a convolutional
 579 torso followed by an LSTM [32] and a dueling head [62]. When leveraging the behavior of
 580 the pre-trained policy to solve new tasks, we instantiate a new network with independent weights
 581 (c.f. Figure 7). One can initialize some of the components of the new network using pre-trained
 582 weights without tying their values (as in common fine-tuning approaches).

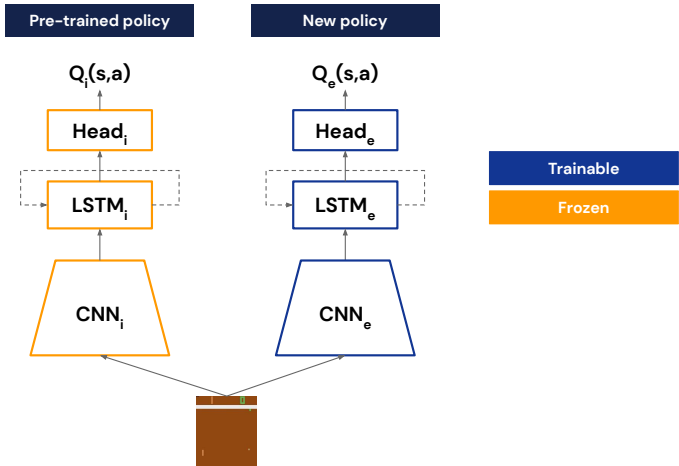


Figure 7: Q-Network architecture for the reinforcement learning stage. The networks use independent sets of parameters, and the weights of the pre-trained policy are kept fixed to preserve the learned behavior.

583 Table 1 summarizes the main hyperparameters of our method. The pre-trained policies were optimized
 584 using Retrace [45]. Learning with rewards was performed with Peng’s $Q(\lambda)$ [50] instead, which we
 585 found to be much more data efficient in our experiments. The reason for this difference is that the
 586 benefits of $Q(\lambda)$ were observed once unsupervised policies had been trained on all Atari games.

Table 1: Hyperparameter values used in R2D2-based agents. The rest of hyperparameters use the values reported by Kapturowski et al. [37].

Hyperparameter	Value
Number of actors	256
Actor parameter update interval	400 environment steps
Sequence length	160 (without burn-in)
Replay buffer size	12.5×10^4 part-overlapping sequences
Priority exponent	0.9
Importance sampling exponent	0
Learning rule (downstream tasks)	$Q(\lambda)$, $\lambda = 0.7$
Learning rule (NGU pre-training)	Retrace(λ), $\lambda = 0.95$
Discount (downstream tasks)	0.99
Discount (NGU pre-training)	0.99
Minibatch size	64
Optimizer	Adam
Optimizer settings	$\epsilon = 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$
Learning rate	2×10^{-4}
Target network update interval	1500 updates
ϵ_{levy} distribution	Log-Uniform[0, 0.1]
Flight length distribution	Zeta with $\mu = 2$

587 It should be noted that our ϵz -greedy baseline under-performs relative to Dabney et al. [14]. This
588 is due to our hyper-parameters and setting being derived from Puigdomènech Badia et al. [52],
589 which adopts the standard Atari pre-processing (e.g. gray scale images and frame stacking). In
590 contrast, Dabney et al. [14] use color images, no frame stacking, a larger neural network and different
591 hyper-parameters (e.g. smaller replay buffer). Studying if the performance of NGU, RND and BT
592 is preserved in this setting is an important direction for future work. We suspect that improving the
593 performance of our ϵz -greedy ablation will also improve our method, since exploration flights are
594 central to both.

595 **C Extended Unsupervised RL Results**

596 We compare the results of our unsupervised pre-training stage against other unsupervised approaches,
 597 standard RL algorithms in the low-data regime and methods that perform unsupervised pre-training
 598 followed by an adaptation stage. Since the considered intrinsic rewards are non-negative, we consider
 599 a baseline where the agent obtains a constant positive reward at each step in order to measure the
 600 performance of policies that seek to stay alive for as long as possible. Results for this baseline
 601 were already considered by Hansen et al. [28] (*Pos Reward NSQ*), but we run our own version of
 602 this baseline using the distributed setting and longer pre-training of 16B frames considered in our
 603 experiments (*Pos Reward R2D2*). Table 2 shows that unsupervised RND and NGU outperform all
 604 baselines by a large margin, confirming the intuition that exploration is a good pre-training objective
 605 for the Atari benchmark. These results suggest that there is a strong correlation between exploration
 606 and the goals established by game designers [10]. In spite of the strong results, it is worth noting
 607 that unsupervised RND and NGU achieve lower scores than random policies in some games, and
 608 can be quite inefficient at collecting rewards in some environments (e.g. they needs long episodes to
 609 obtain high scores). These observations motivate the development of techniques to leverage these
 610 pre-trained policies without compromising performance even when there exists a misalignment
 611 between objectives.

Table 2: Atari Suite comparisons, adapted from Hansen et al. [28] and Liu and Abbeel [41]. @ N represents the amount of RL interaction with reward utilized, with four frames observed at each iteration. *Mdn* and *M* are median and mean human normalized scores, respectively; > 0 is the number of games with better than random performance; and $> H$ is the number of games with human-level performance as defined in Mnih et al. [43]. **Top**: unsupervised learning only. **Mid**: data-limited RL. **Bottom**: RL with unsupervised pre-training.

Algorithm	26 Game Subset Kaiser et al. [35]				47 Game Subset Burda et al. [10]				Full 57 Games Mnih et al. [43]			
	Mdn	M	>0	$>H$	Mdn	M	>0	$>H$	Mdn	M	>0	$>H$
IDF Curiosity @0	-	-	-	-	8.46	24.51	34	5	-	-	-	-
RF Curiosity @0	-	-	-	-	7.32	29.03	36	6	-	-	-	-
Pos Reward NSQ @0	2.18	50.33	14	5	0.69	57.65	26	8	0.29	41.19	28	8
Pos Reward R2D2 @0	9.44	59.55	21	4	14.16	57.53	39	5	3.46	45.23	46	5
Q-DIAYN-5 @0	0.17	-3.60	13	0	0.33	-1.23	25	2	0.34	-2.18	30	2
Q-DIAYN-50 @0	-1.65	-21.77	4	0	-1.69	-16.26	8	0	-3.16	-20.31	9	0
VISR @0	5.60	81.65	19	5	4.04	58.47	35	7	3.77	49.66	40	7
RND@0	48.35	334.65	23	8	41.28	259.43	40	14	40.86	243.01	47	16
NGU @0	80.92	494.54	25	12	96.10	310.27	45	23	81.72	320.06	52	27
SimPLe @100k	9.79	36.20	26	4	-	-	-	-	-	-	-	-
DQN @10M	27.80	52.95	25	7	9.91	28.07	41	7	8.61	27.55	48	7
DQN @200M	100.76	267.51	26	13	-	-	-	-	80.81	239.29	46	20
Rainbow @100k	2.23	10.12	25	1	-	-	-	-	-	-	-	-
PPO @500k	20.93	43.74	25	7	-	-	-	-	-	-	-	-
NSQ @10M	8.20	33.80	22	3	7.29	29.47	37	4	6.80	28.51	43	5
SPR @100k	41.50	70.40	-	7	-	-	-	-	-	-	-	-
CURL @100k	17.50	38.10	-	2	-	-	-	-	-	-	-	-
DrQ @100k	28.42	35.70	-	2	-	-	-	-	-	-	-	-
Q-DIAYN-5 @100k	0.01	16.94	13	2	1.31	19.64	28	6	1.55	16.65	33	6
Q-DIAYN-50 @100k	-1.64	-27.88	3	0	-1.66	-16.74	8	0	-2.53	-24.13	9	0
RF VISR @100k	7.24	58.23	20	6	3.81	42.60	33	9	2.16	35.29	39	9
VISR @100k	9.50	128.07	21	7	9.42	121.08	35	11	6.81	102.31	40	11
GPI RF VISR @100k	5.55	58.77	20	5	4.24	48.38	34	9	3.60	40.01	40	10
GPI VISR @100k	6.59	111.23	22	7	11.70	129.76	38	12	8.99	109.16	44	12
MEPOL @100k	0.34	17.94	-	2	-	-	-	-	-	-	-	-
APT @100k	47.50	69.55	-	7	-	-	-	-	33.41	47.73	-	12

612 **D Extended Atari-57 Results**

Table 3: Atari Suite comparisons for R2D2-based agents. @ N represents the amount of frames with reward utilized, with four frames observed per RL interaction. Mdn , M and CM are median, mean and mean capped human normalized scores, respectively.

Algorithm	Full 57 Games			Hard Exploration		
	Mdn	M	CM	Mdn	M	CM
R2D2 @1B	229.75	864.69	84.56	31.07	39.40	34.75
R2D2 + ϵz -greedy @1B	204.52	578.73	85.11	42.55	53.90	46.21
R2D2 + BT(π_{NGU}) @1B	273.49	1517.13	86.38	100.89	94.20	63.95
R2D2 + BT(π_{RND}) @1B	280.04	1396.78	87.43	93.52	86.75	67.40
R2D2 @5B	490.12	1742.92	90.37	32.49	67.41	44.74
R2D2 + ϵz -greedy @5B	418.41	1275.86	92.49	103.62	95.46	67.85
R2D2 + BT(π_{NGU}) @5B	538.50	2262.21	93.31	193.15	160.02	76.92
R2D2 + BT(π_{RND}) @5B	571.57	2304.19	92.03	144.78	123.38	76.93

Table 4: Atari Suite comparisons with rewards for R2D2-based agents with different amounts of transfer via weights at 5B training frames. Policies are composed of a CNN encoder followed by an LSTM and a dueling head. We compare training from scratch, loading all weights (Full π_{NGU} init) or all weights except those in the dueling head (Partial π_{NGU} init). Mdn , M and CM are median, mean and mean capped human normalized scores, respectively. (**Top**) Without BT. (**Bottom**) With BT(π_{NGU}).

Algorithm	Full 57 Games			Hard Exploration		
	Mdn	M	CM	Mdn	M	CM
R2D2, from scratch	490.12	1742.92	90.37	32.49	67.41	44.74
R2D2, partial π_{NGU} init	668.80	2020.81	93.00	109.33	123.40	67.18
R2D2, full π_{NGU} init	507.58	2359.25	89.91	104.98	101.52	66.20
R2D2 + BT(π_{NGU}), from scratch	538.50	2262.21	93.31	193.15	160.02	76.92
R2D2 + BT(π_{NGU}), partial π_{NGU} init	626.34	1966.83	94.07	200.32	164.54	76.93
R2D2 + BT(π_{NGU}), full π_{NGU} init	529.78	2467.02	92.79	168.18	137.65	76.93

Table 5: Human normalized scores after 5B frames with rewards for R2D2-based agents at different percentiles. Note that the 50th percentile corresponds to the median score across the 57 games. We compare training from scratch, loading all weights (Full π_{NGU} init) or all weights except those in the dueling head (Partial π_{NGU} init).

Method	Percentile				
	50th	40th	20th	10th	5th
R2D2, from scratch	490.12	220.97	132.77	92.31	25.57
R2D2 + BT(π_{NGU}), from scratch	538.50	316.30	163.20	104.41	65.17
R2D2 + BT(π_{RND}), from scratch	571.57	279.97	133.61	87.05	54.27
R2D2, partial π_{NGU} init	668.80	434.02	164.62	105.08	53.59
R2D2 + BT(π_{NGU}), partial π_{NGU} init	626.34	489.14	171.69	113.92	93.52
R2D2, full π_{NGU} init	507.58	293.54	137.79	59.72	41.06
R2D2 + BT(π_{NGU}), full π_{NGU} init	529.78	384.39	163.55	123.56	51.98

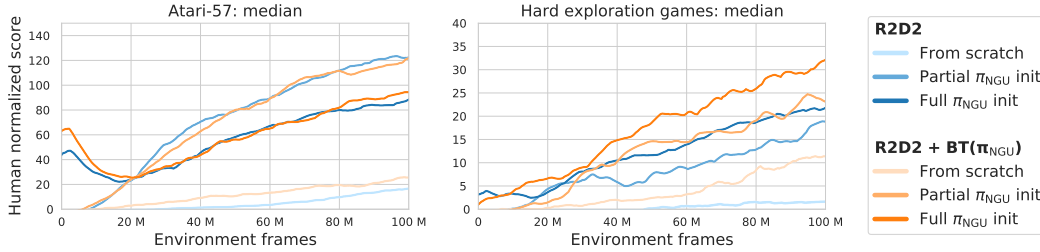


Figure 8: Median human normalized scores for R2D2-based agents with different amounts of transfer via weights during the first 100M frames of training. Policies are composed of a CNN encoder followed by an LSTM and a dueling head. We compare training from scratch, loading the CNN and the LSTM (Partial π_{NGU} init), and loading all weights including the dueling head (Full π_{NGU} init). **(Left)** Full Atari suite. **(Right)** Subset of hard exploration games.

613 E Alternative Reward Functions

614 MsPacman: eating ghosts

- 615 • Pac-dots: 0 points (easy) or -10 points (hard)
- 616 • Eating vulnerable ghosts:
 - 617 – #1 in succession: 200 points
 - 618 – #2 in succession: 400 points
 - 619 – #3 in succession: 800 points
 - 620 – #4 in succession: 1600 points
- 621 • Other actions: 0 points

622 Hero: rescuing miners

- 623 • Dynamiting walls: 0 points (easy) or -300 points (hard)
- 624 • Rescuing a miner: 1000 points
- 625 • Other actions: 0 points

626 F Distributed setting

627 All experiments are run using a distributed setting. The evaluation we do is also identical to the
 628 one done in R2D2 [37]: parallel evaluation workers, which share weights with actors and learners,
 629 run the Q-network against the environment. This worker and all the actor workers are the two
 630 types of workers that draw samples from the environment. For Atari, we apply the standard DQN
 631 pre-processing, as used in R2D2. The next subsections describe how actors, evaluators, and learner
 632 are run in each stage.

633 F.1 Unsupervised stage

634 The computation of the intrinsic NGU reward, r_t^{NGU} , follows the method described in Puig-
 635 domènech Badia et al. [52, Appendix A.1]. In particular, we use the version that combines episodic
 636 intrinsic rewards with the intrinsic reward from Random Network Distillation (RND) [11].

637 We now describe the distributed setup used for NGU, which is largely the same as the one used for
 638 RND. Note that RND can be recovered by removing the components needed for the episodic reward.

639 Learner

- 640 • Sample from the replay buffer a sequence of intrinsic rewards r_t^{NGU} , observations x and
 641 actions a .
- 642 • Use Q-network to learn from (r_t^{NGU}, x, a) with Retrace [45] using the same procedure as in
 643 R2D2.

- 644 • Use last 5 frames of the sampled sequences to train the action prediction network in NGU.
- 645 This means that, for every batch of sequences, all time steps are used to train the RL loss,
- 646 whereas only 5 time steps per sequence are used to optimize the action prediction loss.
- 647 • Use last 5 frames of the sampled sequences to train the predictor of RND.

648 Actor

- 649 • Obtain x_t and r_{t-1}^{NGU} .
- 650 • With these inputs, compute forward pass of R2D2 to obtain a_t .
- 651 • With x_t , compute r_t^{NGU} using the embedding network in NGU.
- 652 • Insert x_t , a_t and r_t^{NGU} in the replay buffer.
- 653 • Step on the environment with a_t .

654 Evaluator

- 655 • Obtain x_t and r_{t-1}^{NGU} .
- 656 • With these inputs, compute forward pass of R2D2 to obtain a_t .
- 657 • With x_t , compute r_t^{NGU} using the embedding network in NGU.
- 658 • Step on the environment with a_t .

659 Distributed training

660 As in R2D2, we train the agent with a single GPU-based learner and a fixed discount factor γ . All
 661 actors collect experience using the same policy, but with a different value of ϵ . This differs from the
 662 original NGU agent, where each actor runs a policy with a different degree of exploratory behavior
 663 and discount factor.

664 In the replay buffer, we store fixed-length sequences of (x, a, r) tuples. These sequences never cross
 665 episode boundaries. Given a single batch of trajectories we unroll both online and target networks on
 666 the same sequence of states to generate value estimates. We use prioritized experience replay with
 667 the same prioritization scheme proposed in [37].

668 F.2 Transfer with BT

669 Learner

- 670 • Sample from the replay buffer a sequence of extrinsic rewards r_t , observations x and actions
 671 a .
- 672 • (expanded action set) Duplicate transitions collected with π_p and relabel the duplicates with
 673 the primitive action taken by π_p when acting.
- 674 • Use Q-network to learn from (r_t, x, a) with Peng’s $Q(\lambda)$ [50] using the same procedure as
 675 in R2D2.

676 Actor

- 677 • (once per episode) Sample ϵ_{levy} .
- 678 • Obtain x_t .
- 679 • If not on a flight, start one with probability ϵ_{levy} .
- 680 • If on a flight, compute forward pass with π_p to obtain a_t . Otherwise, compute forward pass
 681 of R2D2 to obtain a_t . If $a_t = |\mathcal{A}| + 1$, $a_t \leftarrow \pi_p(x)$.
- 682 • Insert x_t , a_t and r_t in the replay buffer.
- 683 • Step on the environment with a_t .

684 **Evaluator**

- 685 • Obtain x_t .
- 686 • Compute forward pass of R2D2 to obtain a_t . If $a_t = |\mathcal{A}| + 1$, $a_t \leftarrow \pi_p(x)$.
- 687 • Step on the environment with a_t .

688 **Distributed training**

689 As in R2D2, we train the agent with a single GPU-based learner and a fixed discount factor γ . All
690 actors collect experience using the same policy, but with a different value of ϵ .

691 In the replay buffer, we store fixed-length sequences of (x, a, r) tuples. These sequences never cross
692 episode boundaries. Given a single batch of trajectories we unroll both online and target networks on
693 the same sequence of states to generate value estimates. We use prioritized experience replay with
694 the same prioritization scheme proposed in [37].

695 **G Intrinsic Rewards**

696 **G.1 Random Network Distillation**

697 The RND [11] intrinsic reward is computed by introducing a random, untrained convolutional network
698 $g : \mathcal{S} \rightarrow \mathbb{R}^d$, and training a network $\hat{g} : \mathcal{S} \rightarrow \mathbb{R}^d$ to predict the outputs of g on all the observations
699 that are seen during training by minimizing the prediction error $\text{err}_{\text{RND}}(s_t) = \|\hat{g}(s_t; \theta) - g(s_t)\|^2$
700 with respect to θ . The intuition is that the prediction error will be large on states that have been visited
701 less frequently by the agent. The dimensionality of the random embedding, d , is a hyperparameter of
702 the algorithm.

703 The RND intrinsic reward is obtained by normalising the prediction error. In this work, we use a
704 slightly different normalization from that reported in [11]. The RND reward at time t is given by

$$r_t^{\text{RND}} = \frac{\text{err}_{\text{RND}}(s_t)}{\sigma_e} \tag{2}$$

705 where σ_e is the running standard deviation of $\text{err}_{\text{RND}}(s_t)$.

706 **G.2 Never Give Up**

707 The NGU intrinsic reward modulates an episodic intrinsic reward, r_t^{episodic} , with a life long signal α_t :

$$r_t^{\text{NGU}} = r_t^{\text{episodic}} \cdot \min \{ \max \{ \alpha_t, 1 \}, L \}, \tag{3}$$

708 where L is a fixed maximum reward scaling. The life-long novelty signal is computed using RND
709 with the normalisation:

$$\alpha_t = \frac{\text{err}_{\text{RND}}(s_t) - \mu_e}{\sigma_e} \tag{4}$$

710 where $\text{err}_{\text{RND}}(x_t)$ is the prediction error described in Appendix G.1, and μ_e and σ_e are its running
711 mean and standard deviation, respectively. The episodic intrinsic reward at time t is computed
712 according to formula:

$$r_t^{\text{episodic}} = \frac{1}{\sqrt{\sum_{f(s_i) \in N_k} K(f(s_t), f(s_i)) + c}} \tag{5}$$

713 where N_k is the set containing the k -nearest neighbors of $f(s_t)$ in M , c is a constant and $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^+$
714 is a kernel function satisfying $K(x, x) = 1$ (which can be thought of as approximating
715 pseudo-counts [52]). Algorithm 4 shows a detailed description of how the episodic intrinsic reward is
716 computed. Below we describe the different components used in Algorithm 4:

- 717 • M : episodic memory containing at time t the previous embeddings
718 $\{f(s_0), f(s_1), \dots, f(s_{t-1})\}$. This memory starts empty at each episode
- 719 • k : number of nearest neighbours

- 720 • $N_k = \{f(s_i)\}_{i=1}^k$: set of k -nearest neighbours of $f(s_t)$ in the memory M ; we call $N_k[i] =$
- 721 $f(s_i) \in N_k$ for ease of notation
- 722 • K : kernel defined as $K(x, y) = \frac{\epsilon}{\frac{d^2(x, y)}{d_m^2} + \epsilon}$ where ϵ is a small constant, d is the Euclidean
- 723 distance and d_m^2 is a running average of the squared Euclidean distance of the k -nearest
- 724 neighbors
- 725 • c : pseudo-counts constant
- 726 • ξ : cluster distance
- 727 • s_m : maximum similarity

Algorithm 4: Computation of the episodic intrinsic reward at time t : r_t^{episodic} .

Input : M ; k ; $f(s_t)$; c ; ϵ ; ξ ; s_m ; d_m^2

Output : r_t^{episodic}

Compute the k -nearest neighbours of $f(s_t)$ in M and store them in a list N_k

Create a list of floats d_k of size k

/* The list d_k will contain the distances between the embedding $f(s_t)$ and its neighbours N_k . */

for $i \in \{1, \dots, k\}$ **do**

$d_k[i] \leftarrow d^2(f(s_t), N_k[i])$

end

Update the moving average d_m^2 with the list of distances d_k

/* Normalize the distances d_k with the updated moving average d_m^2 . */

$d_n \leftarrow \frac{d_k}{d_m^2}$

/* Cluster the normalized distances d_n i.e. they become 0 if too small and 0_k is a list of k zeros. */

$d_n \leftarrow \max(d_n - \xi, 0_k)$

/* Compute the Kernel values between the embedding $f(s_t)$ and its neighbours N_k . */

$K_v \leftarrow \frac{\epsilon}{d_n + \epsilon}$

/* Compute the similarity between the embedding $f(s_t)$ and its neighbours N_k . */

$s \leftarrow \sqrt{\sum_{i=1}^k K_v[i] + c}$

/* Compute the episodic intrinsic reward at time t : r_t^i . */

if $s > s_m$ **then**

$r_t^{\text{episodic}} \leftarrow 0$

else

$r_t^{\text{episodic}} \leftarrow 1/s$

Table 6: Results per game for R2D2-based agents at 5B training frames.

Game	R2D2	R2D2 + ϵ -greedy	R2D2 + BT(π_{NGU})	R2D2 + BT(π_{RND})
alien	10831.17 ± 2114.29	14634.02 ± 1109.15	15657.57 ± 1717.96	12844.24 ± 1447.72
amidar	11761.67 ± 1560.86	6784.28 ± 718.05	10394.96 ± 891.60	7730.43 ± 670.76
assault	15940.72 ± 3531.69	9177.28 ± 2170.26	15060.31 ± 740.63	11533.24 ± 809.49
asterix	472812.21 ± 222663.81	374966.62 ± 135810.51	630663.91 ± 82753.46	468724.08 ± 120822.86
asteroids	45716.28 ± 3642.38	147005.85 ± 44313.45	31957.42 ± 15540.09	37455.64 ± 9263.04
atlantis	1514724.43 ± 10941.36	1132188.04 ± 43551.36	1491384.23 ± 5978.05	1545954.35 ± 9001.60
bank heist	965.63 ± 133.72	1058.75 ± 135.46	13913.32 ± 3529.15	82132.27 ± 101709.64
battle zone	292553.41 ± 18196.77	312367.76 ± 43554.18	258533.57 ± 22865.64	285925.87 ± 44912.86
beam rider	18472.45 ± 1977.78	22403.95 ± 1596.92	16301.02 ± 1853.73	15619.99 ± 2048.77
berzerk	12343.83 ± 3331.54	3846.56 ± 1723.24	8359.80 ± 201.10	14687.68 ± 401.76
bowling	141.64 ± 4.52	156.32 ± 8.11	174.27 ± 0.10	196.01 ± 57.42
boxing	99.96 ± 0.03	99.94 ± 0.06	100.00 ± 0.00	99.98 ± 0.03
breakout	432.65 ± 27.35	393.19 ± 35.12	441.21 ± 15.08	429.38 ± 15.52
centipede	189502.66 ± 31388.08	358841.20 ± 73578.20	178635.17 ± 17227.15	196880.46 ± 24278.18
chopper command	611393.11 ± 65206.69	697655.53 ± 215090.74	573055.88 ± 75343.57	797052.58 ± 52012.04
crazy climber	229992.57 ± 17738.33	212001.76 ± 1853.07	226821.26 ± 3608.19	198736.17 ± 7631.83
defender	547238.15 ± 2579.38	516521.06 ± 11969.59	540124.74 ± 4488.40	524003.44 ± 1316.59
demon attack	143662.42 ± 88.16	141352.18 ± 3848.73	143762.91 ± 106.75	143578.47 ± 25.05
double dunk	23.99 ± 0.02	23.88 ± 0.06	23.85 ± 0.15	23.93 ± 0.05
enduro	2358.37 ± 3.32	2359.08 ± 1.03	2361.56 ± 1.03	2350.39 ± 8.42
fishing derby	12.80 ± 77.79	64.74 ± 0.59	52.58 ± 0.32	62.11 ± 5.59
freeway	33.87 ± 0.08	33.77 ± 0.03	33.79 ± 0.08	33.79 ± 0.07
frostbite	9287.24 ± 167.11	8504.41 ± 940.72	17692.42 ± 2871.83	9419.45 ± 188.92
gopher	117398.58 ± 2485.82	84140.40 ± 12919.83	113716.78 ± 3966.91	94670.35 ± 2285.63
gravitar	6123.08 ± 103.19	5798.68 ± 735.59	8373.70 ± 1260.75	7428.57 ± 2459.91
hero	46048.07 ± 6970.26	39700.22 ± 4379.84	40825.09 ± 3736.25	42959.86 ± 7950.56
ice hockey	32.43 ± 30.64	30.65 ± 28.17	60.36 ± 4.94	57.96 ± 0.90
jamesbond	6056.14 ± 1643.52	3843.92 ± 118.35	1484.87 ± 489.66	2870.03 ± 907.76
kangaroo	14672.37 ± 187.16	14730.99 ± 114.20	15965.79 ± 36.61	15128.66 ± 188.17
krull	10081.04 ± 594.10	10171.52 ± 399.81	406596.00 ± 55547.76	316960.78 ± 217091.10
kung fu master	200721.64 ± 2265.35	171591.29 ± 8516.87	196638.89 ± 456.09	610699.23 ± 60053.99
montezuma revenge	1478.38 ± 1114.20	1467.77 ± 1104.72	12086.71 ± 1217.76	6266.67 ± 471.40
ms pacman	11212.85 ± 103.23	7511.39 ± 406.77	10996.90 ± 262.74	10656.00 ± 356.46
name this game	32138.12 ± 2156.95	37343.04 ± 1917.73	30252.11 ± 884.84	28746.14 ± 1798.77
phoenix	712101.72 ± 62738.09	80611.18 ± 25316.56	553429.34 ± 24278.55	283686.99 ± 172323.63
pitfall	-0.19 ± 0.15	-12.34 ± 4.20	-0.39 ± 0.39	-0.03 ± 0.04
pong	20.93 ± 0.01	20.49 ± 0.10	20.90 ± 0.01	20.94 ± 0.01
private eye	23592.22 ± 11876.55	50770.82 ± 14984.92	40435.54 ± 51.04	40480.67 ± 38.23
qbert	24343.75 ± 1904.89	16975.13 ± 1332.44	16057.31 ± 318.87	10990.08 ± 7241.50
riverraid	32325.07 ± 1185.15	30582.53 ± 638.47	28550.32 ± 2298.03	30566.86 ± 1764.50
road runner	423191.07 ± 53071.15	88890.04 ± 24971.18	251261.09 ± 31741.38	248661.22 ± 19416.63
robotank	97.23 ± 1.22	108.92 ± 4.79	98.45 ± 2.85	100.57 ± 6.32
seaquest	188771.84 ± 20759.57	175745.09 ± 120718.82	86605.86 ± 55065.85	38185.98 ± 22949.18
skiing	-29854.11 ± 85.79	-30060.81 ± 142.32	-30121.95 ± 70.62	-29589.38 ± 69.40
solaris	17741.02 ± 5340.46	16127.73 ± 2975.20	24366.59 ± 4868.05	18727.45 ± 4806.17
space invaders	3621.76 ± 5.81	3547.78 ± 35.31	30609.21 ± 7141.11	46704.49 ± 7017.79
star gunner	223536.63 ± 48548.34	179698.69 ± 12194.36	171294.31 ± 23185.79	156691.39 ± 16704.19
surround	8.24 ± 0.48	1.48 ± 8.12	5.86 ± 1.44	-3.62 ± 4.79
tennis	7.99 ± 22.56	7.98 ± 22.51	23.96 ± 0.01	7.97 ± 22.56
time pilot	139931.67 ± 70521.78	71768.84 ± 2933.22	44936.87 ± 137.49	77711.97 ± 4735.53
tutankham	324.02 ± 4.26	311.65 ± 8.62	420.36 ± 30.13	357.26 ± 14.22
up n down	529363.05 ± 16813.20	394984.70 ± 34313.42	562739.02 ± 8527.59	585355.01 ± 4718.67
venture	0.00 ± 0.00	1833.85 ± 43.73	2110.64 ± 55.39	1910.15 ± 13.98
video pinball	454023.46 ± 377076.03	107071.98 ± 67142.18	463141.28 ± 426927.92	646671.78 ± 403584.41
wizard of wor	40833.65 ± 4776.81	38275.31 ± 4177.41	30453.12 ± 2470.20	30399.63 ± 2345.83
yars revenge	279765.86 ± 27370.20	250483.70 ± 54593.32	280333.48 ± 69704.31	200850.59 ± 72885.67
zaxxon	56059.14 ± 3217.77	66099.28 ± 8520.19	67611.78 ± 6226.04	59926.08 ± 5834.47

Table 7: Results per game for R2D2 agents with different amounts of transfer via weights at 5B training frames. Policies are composed of a CNN encoder followed by an LSTM and a dueling head. We compare training from scratch, loading all weights (Full π_{NGU} init) or all weights except those in the dueling head (Partial π_{NGU} init).

Game	From scratch	Partial π_{NGU} init	Full π_{NGU} init
alien	10831.17 \pm 2114.29	27299.78 \pm 5730.57	18027.35 \pm 6731.75
amidar	11761.67 \pm 1560.86	13647.09 \pm 3380.90	3518.30 \pm 2353.96
assault	15940.72 \pm 3531.69	14653.32 \pm 2047.43	12533.61 \pm 1001.68
asterix	472812.21 \pm 222663.81	789344.47 \pm 80638.00	676662.54 \pm 8536.94
asteroids	45716.28 \pm 3642.38	73298.12 \pm 22688.38	23127.43 \pm 5425.84
atlantis	1514724.43 \pm 10941.36	1537659.81 \pm 7693.86	1556234.51 \pm 9709.74
bank heist	965.63 \pm 133.72	1841.77 \pm 52.75	5816.24 \pm 3137.60
battle zone	292553.41 \pm 18196.77	301715.60 \pm 12875.97	248939.89 \pm 31788.00
beam rider	18472.45 \pm 1977.78	16179.19 \pm 4179.36	11040.66 \pm 799.77
berzerk	12343.83 \pm 3331.54	16888.63 \pm 2330.72	25465.10 \pm 9886.89
bowling	141.64 \pm 4.52	170.36 \pm 16.55	180.78 \pm 2.94
boxing	99.96 \pm 0.03	99.97 \pm 0.05	99.94 \pm 0.07
breakout	432.65 \pm 27.35	520.19 \pm 64.65	487.51 \pm 49.14
centipede	189502.66 \pm 31388.08	528000.27 \pm 10403.62	500534.38 \pm 9267.05
chopper command	611393.11 \pm 65206.69	937637.69 \pm 57836.24	764150.71 \pm 44756.44
crazy climber	229992.57 \pm 17738.33	275735.70 \pm 15244.51	246498.24 \pm 12319.58
defender	547238.15 \pm 2579.38	534656.86 \pm 2880.53	523660.11 \pm 2604.26
demon attack	143662.42 \pm 88.16	143592.16 \pm 77.27	143574.75 \pm 69.35
double dunk	23.99 \pm 0.02	23.99 \pm 0.02	23.83 \pm 0.06
enduro	2358.37 \pm 3.32	2359.39 \pm 8.12	2353.16 \pm 1.30
fishing derby	12.80 \pm 77.79	68.70 \pm 2.46	59.22 \pm 2.55
freeway	33.87 \pm 0.08	33.83 \pm 0.06	33.79 \pm 0.04
frostbite	9287.24 \pm 167.11	161595.33 \pm 32917.44	10307.96 \pm 1087.09
gopher	117398.58 \pm 2485.82	113094.41 \pm 4837.16	102781.75 \pm 9613.77
gravitar	6123.08 \pm 103.19	7090.19 \pm 1359.52	5174.90 \pm 544.76
hero	46048.07 \pm 6970.26	43982.29 \pm 4124.79	40628.07 \pm 4008.99
ice hockey	32.43 \pm 30.64	69.57 \pm 1.18	47.67 \pm 10.59
jamesbond	6056.14 \pm 1643.52	6109.60 \pm 1643.75	3979.12 \pm 1233.92
kangaroo	14672.37 \pm 187.16	14863.32 \pm 259.85	15192.97 \pm 832.48
krull	10081.04 \pm 594.10	11806.49 \pm 580.05	372307.71 \pm 161921.43
kung fu master	200721.64 \pm 2265.35	200305.15 \pm 5711.26	207401.69 \pm 1755.69
montezuma revenge	1478.38 \pm 1114.20	2666.30 \pm 235.18	2500.00 \pm 0.00
ms pacman	11212.85 \pm 103.23	11795.03 \pm 640.73	11509.67 \pm 563.98
name this game	32138.12 \pm 2156.95	33811.87 \pm 2091.30	29242.89 \pm 1113.73
phoenix	712101.72 \pm 62738.09	812093.31 \pm 42328.98	801952.54 \pm 33211.40
pitfall	-0.19 \pm 0.15	-1.43 \pm 1.17	-0.61 \pm 0.60
pong	20.93 \pm 0.01	20.96 \pm 0.01	20.79 \pm 0.13
private eye	23592.22 \pm 11876.55	30345.57 \pm 10971.52	28653.02 \pm 9512.72
qbert	24343.75 \pm 1904.89	40943.28 \pm 16722.72	62018.46 \pm 34865.08
riverraid	32325.07 \pm 1185.15	35995.19 \pm 825.70	35845.18 \pm 3486.49
road runner	423191.07 \pm 53071.15	311557.84 \pm 59675.36	279988.24 \pm 58503.61
robotank	97.23 \pm 1.22	111.78 \pm 4.63	91.93 \pm 2.09
seaquest	188771.84 \pm 20759.57	629817.31 \pm 145648.54	31735.24 \pm 31257.98
skiing	-29854.11 \pm 85.79	-29550.10 \pm 495.22	-29981.62 \pm 564.13
solaris	17741.02 \pm 5340.46	29751.08 \pm 2076.41	22269.53 \pm 6584.51
space invaders	3621.76 \pm 5.81	41357.74 \pm 8968.52	42695.22 \pm 7148.08
star gunner	223536.63 \pm 48548.34	212821.27 \pm 19723.78	129058.91 \pm 11260.80
surround	8.24 \pm 0.48	6.86 \pm 0.27	-3.29 \pm 8.36
tennis	7.99 \pm 22.56	23.93 \pm 0.02	23.74 \pm 0.16
time pilot	139931.67 \pm 70521.78	65101.20 \pm 7622.18	49957.68 \pm 602.83
tutankham	324.02 \pm 4.26	333.37 \pm 10.62	312.19 \pm 4.32
up n down	529363.05 \pm 16813.20	572472.73 \pm 4512.30	595047.70 \pm 1976.45
venture	0.00 \pm 0.00	1930.36 \pm 32.97	1958.13 \pm 48.98
video pinball	454023.46 \pm 377076.03	113036.77 \pm 3633.15	108849.58 \pm 4753.59
wizard of wor	40833.65 \pm 4776.81	46931.25 \pm 1708.52	19100.00 \pm 1930.29
yars revenge	279765.86 \pm 27370.20	284565.01 \pm 29764.09	294877.82 \pm 52551.36
zaxxon	56059.14 \pm 3217.77	77649.92 \pm 15901.03	75850.01 \pm 10805.13

Table 8: Results per game for R2D2+BT(π_{NGU}) agents with different amounts of transfer via weights at 5B training frames. Policies are composed of a CNN encoder followed by an LSTM and a dueling head. We compare training from scratch, loading all weights (Full π_{NGU} init) or all weights except those in the dueling head (Partial π_{NGU} init).

Game	From scratch	Partial π_{NGU} init	Full π_{NGU} init
alien	15657.57 \pm 1717.96	35441.47 \pm 3848.23	32822.74 \pm 3181.51
amidar	10394.96 \pm 891.60	11564.75 \pm 1726.50	8185.74 \pm 346.35
assault	15060.31 \pm 740.63	12617.35 \pm 3267.80	12992.70 \pm 2783.13
asterix	630663.91 \pm 82753.46	731452.01 \pm 71621.96	815753.95 \pm 91022.04
asteroids	31957.42 \pm 15540.09	53916.46 \pm 12925.49	77368.55 \pm 17900.09
atlantis	1491384.23 \pm 5978.05	1512404.85 \pm 10047.26	1544673.15 \pm 5590.54
bank heist	13913.32 \pm 3529.15	11674.48 \pm 1694.59	8565.69 \pm 4070.27
battle zone	258533.57 \pm 22865.64	321572.13 \pm 32083.18	206177.95 \pm 27251.07
beam rider	16301.02 \pm 1853.73	17465.26 \pm 4954.96	21680.37 \pm 5991.66
berzerk	8359.80 \pm 201.10	15824.26 \pm 5556.26	16161.60 \pm 2848.82
bowling	174.27 \pm 0.10	229.04 \pm 6.36	201.86 \pm 25.21
boxing	100.00 \pm 0.00	99.99 \pm 0.01	99.83 \pm 0.12
breakout	441.21 \pm 15.08	469.52 \pm 31.70	474.30 \pm 37.00
centipede	178635.17 \pm 17227.15	362169.27 \pm 43577.84	525652.45 \pm 19649.69
chopper command	573055.88 \pm 75343.57	766193.62 \pm 109233.61	860939.78 \pm 116076.87
crazy climber	226821.26 \pm 3608.19	224084.35 \pm 7322.03	256189.16 \pm 21605.50
defender	540124.74 \pm 4488.40	525077.01 \pm 6084.08	503190.45 \pm 32182.12
demon attack	143762.91 \pm 106.75	143537.16 \pm 139.44	143554.33 \pm 63.61
double dunk	23.85 \pm 0.15	23.90 \pm 0.07	23.81 \pm 0.11
enduro	2361.56 \pm 1.03	2352.77 \pm 4.01	2353.81 \pm 3.50
fishing derby	52.58 \pm 0.32	64.20 \pm 2.52	52.53 \pm 1.24
freeway	33.79 \pm 0.08	33.69 \pm 0.13	33.57 \pm 0.09
frostbite	17692.42 \pm 2871.83	20847.02 \pm 15492.47	19716.28 \pm 13424.59
gopher	113716.78 \pm 3966.91	105370.92 \pm 12883.78	101383.00 \pm 7891.20
gravitar	8373.70 \pm 1260.75	8358.32 \pm 1022.94	6104.39 \pm 1215.73
hero	40825.09 \pm 3736.25	45837.71 \pm 194.61	43076.13 \pm 4119.55
ice hockey	60.36 \pm 4.94	66.97 \pm 0.72	30.15 \pm 4.43
jamesbond	1484.87 \pm 489.66	1137.12 \pm 89.95	4119.56 \pm 490.29
kangaroo	15965.79 \pm 36.61	15862.51 \pm 234.05	15855.24 \pm 224.30
krull	406596.00 \pm 55547.76	154118.68 \pm 179080.83	350784.05 \pm 205164.12
kung fu master	196638.89 \pm 456.09	193105.00 \pm 5378.00	195990.25 \pm 4969.60
montezuma revenge	12086.71 \pm 1217.76	12714.19 \pm 824.60	11472.65 \pm 629.87
ms pacman	10996.90 \pm 262.74	11337.68 \pm 1176.17	10770.08 \pm 1005.72
name this game	30252.11 \pm 884.84	30656.34 \pm 373.97	28103.24 \pm 2023.96
phoenix	553429.34 \pm 24278.55	510548.66 \pm 236677.91	805269.57 \pm 37169.21
pitfall	-0.39 \pm 0.39	-0.44 \pm 0.32	-0.01 \pm 0.02
pong	20.90 \pm 0.01	20.93 \pm 0.02	20.19 \pm 0.93
private eye	40435.54 \pm 51.04	40472.03 \pm 39.10	40448.67 \pm 40.02
qbert	16057.31 \pm 318.87	15983.72 \pm 888.74	17954.73 \pm 302.05
riverraid	28550.32 \pm 2298.03	34591.91 \pm 831.68	34268.13 \pm 149.84
road runner	251261.09 \pm 31741.38	307342.61 \pm 41017.79	308258.62 \pm 84372.17
robotank	98.45 \pm 2.85	103.82 \pm 2.98	90.17 \pm 8.09
seaquest	86605.86 \pm 55065.85	259408.14 \pm 144362.40	88376.19 \pm 105086.08
skiing	-30121.95 \pm 70.62	-29786.38 \pm 401.06	-29878.47 \pm 289.38
solaris	24366.59 \pm 4868.05	24111.78 \pm 2745.89	19355.27 \pm 4102.09
space invaders	30609.21 \pm 7141.11	43675.67 \pm 6763.52	51318.54 \pm 4277.08
star gunner	171294.31 \pm 23185.79	138390.23 \pm 6320.24	135606.16 \pm 8098.51
surround	5.86 \pm 1.44	6.89 \pm 1.03	-5.48 \pm 0.88
tennis	23.96 \pm 0.01	23.97 \pm 0.01	23.86 \pm 0.06
time pilot	44936.87 \pm 137.49	65721.81 \pm 3213.79	53053.12 \pm 4275.77
tutankham	420.36 \pm 30.13	385.05 \pm 5.02	342.17 \pm 4.79
up n down	562739.02 \pm 8527.59	581518.52 \pm 5198.05	582923.55 \pm 2849.01
venture	2110.64 \pm 55.39	2308.26 \pm 14.97	2054.24 \pm 65.41
video pinball	463141.28 \pm 426927.92	133315.36 \pm 70576.86	640269.38 \pm 330619.65
wizard of woe	30453.12 \pm 2470.20	34648.51 \pm 4182.32	26076.64 \pm 2060.62
yars revenge	280333.48 \pm 69704.31	320777.97 \pm 64750.83	267861.48 \pm 81193.44
zaxxon	67611.78 \pm 6226.04	75165.80 \pm 4030.42	82868.90 \pm 10160.99

Table 9: Final scores per game in our ablation study after 5B frames. We consider versions of $\text{BT}(\pi_{\text{NGU}})$ where the pre-trained policy is used for temporally-extended exploitation (flights), as an extra action (action), or both.

Game	R2D2	R2D2 + $\text{BT}(\pi_{\text{NGU}})$ (flights)	R2D2 + $\text{BT}(\pi_{\text{NGU}})$ (action)	R2D2 + $\text{BT}(\pi_{\text{NGU}})$
asterix	472812.21 \pm 222663.81	630663.91 \pm 82753.46	512869.97 \pm 109039.77	618352.67 \pm 103940.46
bank heist	965.63 \pm 133.72	13913.32 \pm 3529.15	11052.82 \pm 6848.39	12424.39 \pm 1443.95
frostbite	9287.24 \pm 167.11	9114.77 \pm 511.31	10506.15 \pm 3653.44	17692.42 \pm 2871.83
gravitar	6123.08 \pm 103.19	6308.62 \pm 78.84	7228.15 \pm 1842.07	8373.70 \pm 1260.75
jamesbond	6056.14 \pm 1643.52	1615.21 \pm 602.84	3962.10 \pm 798.68	1484.87 \pm 489.66
montezuma revenge	1478.38 \pm 1114.20	11152.10 \pm 664.82	6433.33 \pm 372.68	13265.91 \pm 372.02
ms pacman	11212.85 \pm 103.23	10996.90 \pm 262.74	10648.85 \pm 796.10	10839.60 \pm 445.00
pong	20.93 \pm 0.01	20.90 \pm 0.01	20.94 \pm 0.04	20.88 \pm 0.02
private eye	23592.22 \pm 11876.55	40492.85 \pm 27.35	37029.10 \pm 2437.54	40435.54 \pm 51.04
space invaders	3621.76 \pm 5.81	30671.58 \pm 4418.89	3597.41 \pm 21.06	30609.21 \pm 7141.11
tennis	7.99 \pm 22.56	8.00 \pm 22.54	-7.15 \pm 22.00	23.96 \pm 0.01
up n down	529363.05 \pm 16813.20	562739.02 \pm 8527.59	550665.39 \pm 1852.05	566938.61 \pm 2428.52

Table 10: Final scores per task in Atari games with modified reward functions. We report training results for the standard game reward, a variant with sparse rewards (easy), and a task with deceptive rewards (hard). Despite the pre-trained policy might obtain low or even negative scores in some of the tasks, committing to its exploratory behavior eventually lets the agent discover strategies that lead to high returns.

Game	R2D2	R2D2 + <i>ez-greedy</i>	Fine-tuning π_{NGU}	π_{NGU}	R2D2 + $\text{BT}(\pi_{\text{NGU}})$
Ms Pacman: original	11407 \pm 122	8099 \pm 868	8359 \pm 2117	1360	10984 \pm 665
Ms Pacman: ghosts (easy)	8375 \pm 577	4322 \pm 932	8356 \pm 551	146	8789 \pm 651
Ms Pacman: ghosts (hard)	2836 \pm 26	4018 \pm 1025	1891 \pm 1342	-898	7868 \pm 1085
Hero: original	43762 \pm 4918	39018 \pm 3262	46848 \pm 1199	9298	42675 \pm 3905
Hero: miners (easy)	3000 \pm 0	3000 \pm 0	3000 \pm 0	1351	4665 \pm 470
Hero: miners (hard)	2677 \pm 23	2155 \pm 95	700 \pm 0	-1473	3547 \pm 122

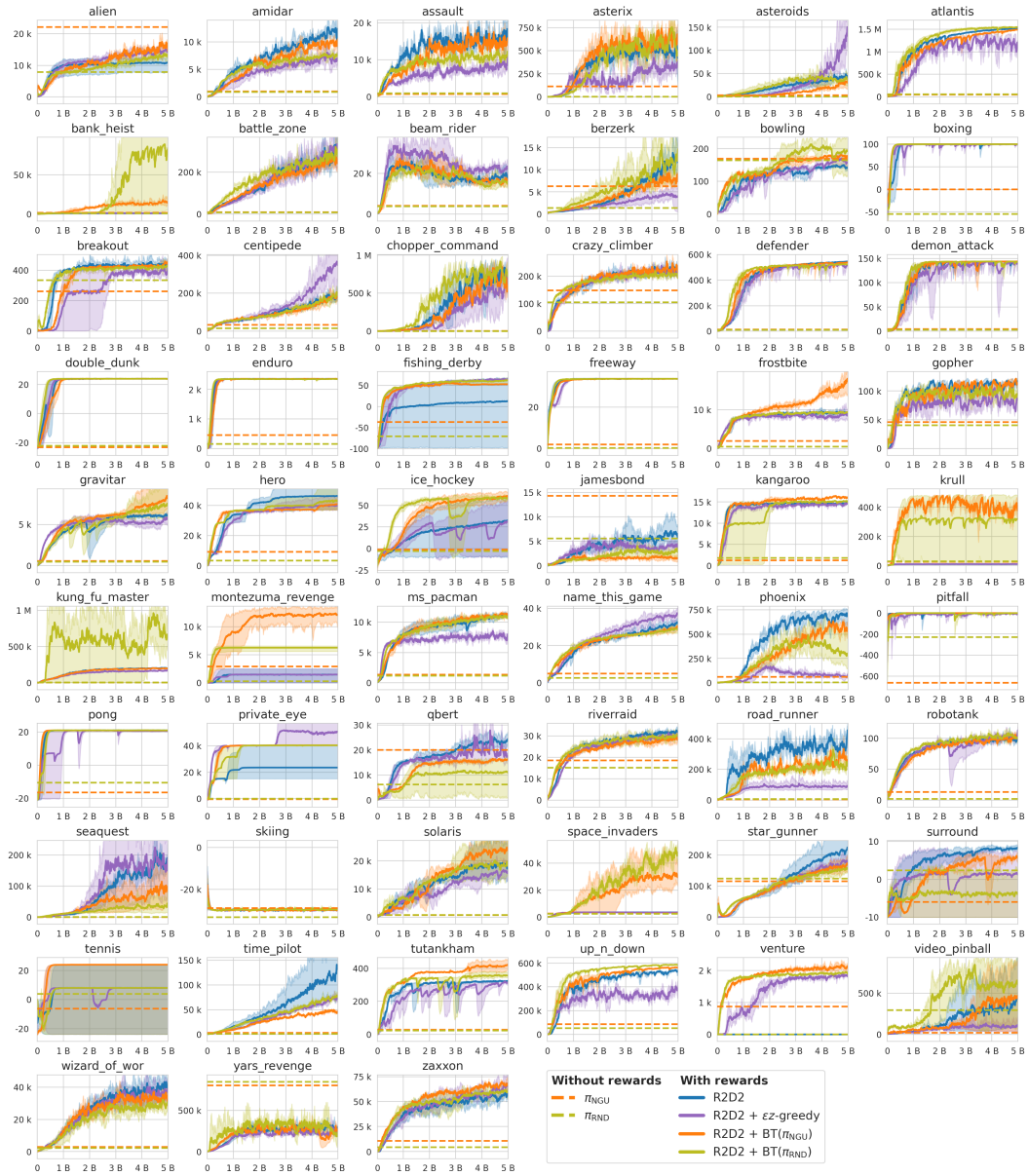


Figure 9: Training curves in all 57 Atari games for R2D2-based agents. Shading shows maximum and minimum over 3 runs, while dark lines indicate the mean.

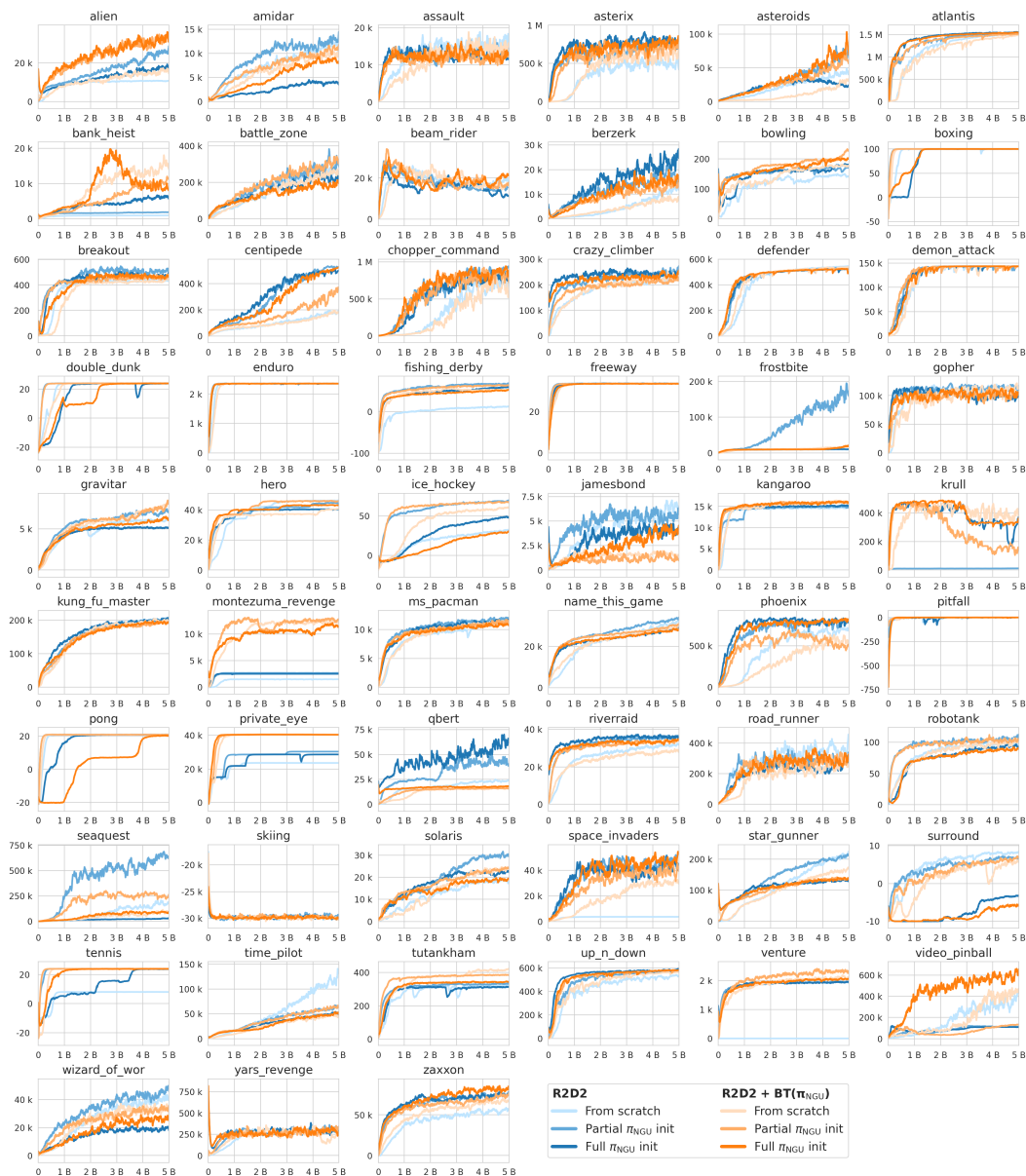


Figure 10: Training curves in all 57 Atari games for R2D2-based agents with different amounts of transfer via weights. Policies are composed of a CNN encoder followed by an LSTM and a dueling head. We compare training from scratch, loading all weights (Full π_{NGU} init) or all weights except those in the dueling head (Partial π_{NGU} init). Shading with maximum and minimum over runs is not shown for clarity, but all plots report the mean over 3 seeds.

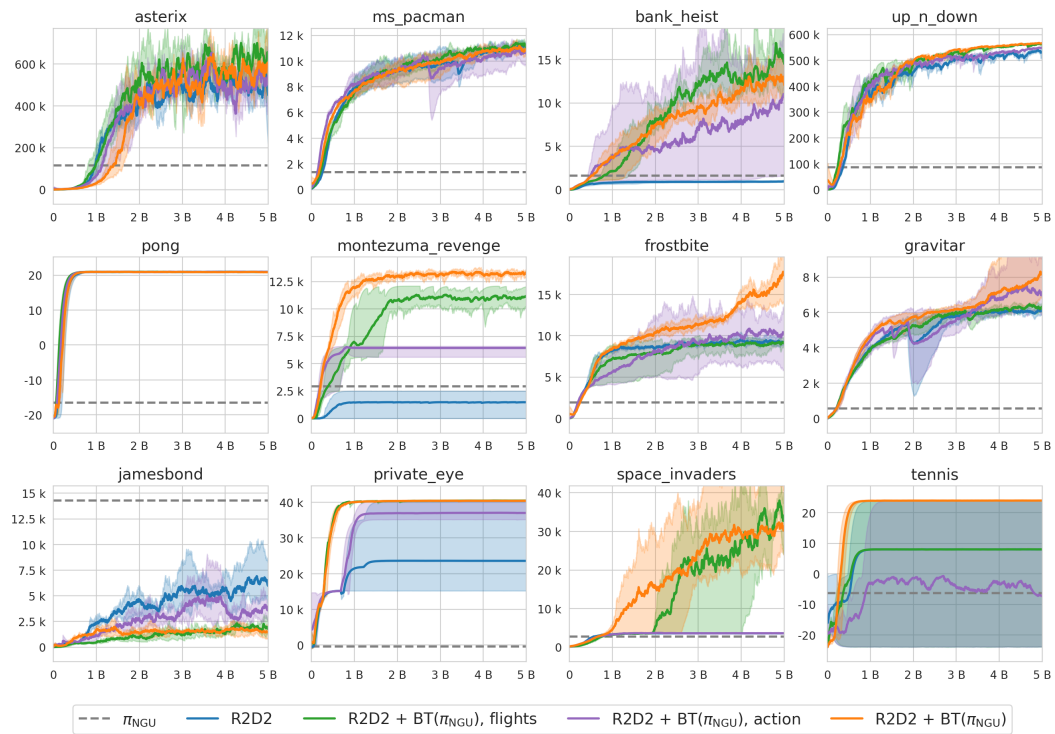


Figure 11: Training curves for ablation experiments. Shading shows maximum and minimum over 3 runs, while dark lines indicate the mean. Both ablations of BT offer benefits over the baseline, but in different sets of games. Combining them retains the best of both methods, and boosts performance even further in some games.