

“Y’all got any more of those Rewards?”

Investigating Transfer Learning in Modern Deep Reinforcement Learning Architectures

Gerrit Bartels
School of Human Sciences
Osnabrück University
Osnabrück, Germany
gbartels@uni-osnabrueck.de

Thorsten Krause
School of Business Administration and Economics
Osnabrück University
Osnabrück, Germany
thkrause@uni-osnabrueck.de

Jacob Dudek
School of Human Sciences
Osnabrück University
Osnabrück, Germany
jadudek@uni-osnabrueck.de

Abstract—Abstract here

Index Terms—reinforcement learning, double deep Q-networks, soft Q-learning, prioritized replay buffer, transfer learning, #publishmoreunsuccessfulresearch, wevehadwaytoolittletimeforthis

I. INTRODUCTION

Through reinforcement learning, machine-based agents autonomously learn how to perform one or multiple tasks at hand. However, traditional reinforcement learning approaches only consider a single deterministic optimal path to achieving the task. Concerned about resulting noise and imperfect behavior, [1] propose a maximum entropy-based, stochastic alternative. The authors in [2] point out additional benefits of maximum entropy-based policies as robustness or exploration for multi-modal objectives. They explain

“Instead of learning the best way to perform the task, [maximum entropy-based] policies try to learn all of the ways of performing the task.” [2]

As existing solutions could only represent a limited range of distributions, [2] further propose *soft Q-learning* as a more general approach capable of asymptotically approximating arbitrary distributions. The authors explicitly developed soft Q-learning for pre-training agents on general tasks and fine-tuning them for more specific tasks later, i.e. for a transfer learning scenario. In reinforcement learning, transfer learning means training an agent on source tasks and leveraging the generated knowledge to perform or learn a target task. For example, [3] later applied soft Q-learning to combine policies trained on different stacking tasks for a robotic agent. [4] investigate performance bounds of soft Q-learning agents in transfer learning scenarios.

However, to the best of our knowledge, researchers have not yet explicitly investigated the general potential advantage of soft Q-learning for transfer learning over deterministic methods. To deepen our understanding on how soft Q-learning compares to deterministic approaches on transfer learning tasks, we seek to answer the following research question:

RQ 1. How does the benefit through transfer learning compare between soft-Q learning and traditional approaches?

During our research, we noticed that models that spent more time on the source task tended to learn the target task more slowly. Researchers have claimed a link between overfitting on the source task and low-performance on the target task *negative transfer* [5], i.e. low-performance on the target task [6]. As we could not detect definite proof for this relationship, we seek to investigate its presence regarding both traditional approaches and soft Q-learning. Knowing that a policy should not train on the source task for too long could save resources for fine-tuning on the target task. Hence, we define a second research question:

RQ 2. Does overfitting on the source task lead to negative transfer in soft Q-learning or in traditional approaches?

We pursued both research questions experimentally on the popular NES game *Super Mario Bros.*. A DDQN and a Soft Q-Network (SQN) agent learned to play one level and reused the final weights to learn to play the second level. Afterwards, we compared their performance to that of agents who only trained on the second level. We are, to the best of our knowledge, the first to directly compare the transfer learning capabilities of soft-Q learning and DDQNs. As both levels are similar, but not identical in both visuals and tasks, they are similar to real world use-cases. Hence, we assume that our findings are applicable to general transfer learning scenarios. With our research, we shed light on soft Q-learning’s transfer learning capabilities and help practitioners take more informed choices between deterministic and stochastic policies for transfer learning. We also give founded advice on how far a policy should train on the source tasks to allow practitioners save expensive resources.

II. RELATED WORK

A. Transfer Learning in Reinforcement Learning

In the reinforcement learning domain, transfer learning allows to leverage knowledge gained in one context to improve

training in another context. For example, it can aim at obtaining more useful initial samples to get a head start in training (*jumpstart improvement*), requiring less samples for training overall (*learning speed improvement*) or final performance (*asymptotic improvement*) [7], [8].

Based on [8] and [9], we shorten the definition of transfer learning in the reinforcement learning context by [10] to

Given a set of source domains \mathcal{M}_s and a target domain \mathcal{M}_t , transfer learning aims to learn an optimal policy π^ for the target domain, by leveraging exterior information from \mathcal{M}_s as well as interior information from \mathcal{M}_t .*

Various taxonomies of transfer learning for reinforcement learning exist and differ both regarding scenarios and methods. For example, one author [7] characterizes scenarios by

- *Setting*: Similarity of state and action space as well as reward dynamics between source and target tasks.
- *Knowledge*: Type of transferred information, i.e. samples, representations or pre-trained policies.
- *Objectives*: Purpose behind employing transfer learning, usually related to enhanced training performance.

while [10] employ a total of six dimensions. By accumulating the taxonomies by [7], [8], [11], we derive the following popular methods for transfer learning in reinforcement learning methods:

- *Starting-Point methods*: Pre-train policy on the source domains and continue training on the target domain.
- *Instance transfer*: Reuse samples from the source domains for training on the target domain.
- *Policy transfer*: Pre-train policies on the source domain. Incorporate their knowledge in training on the target domain, i.e. through weighting schemes.
- *Representation Transfer*: Exploit lower-dimensional representations of state, action or reward dynamics to generalize between domains.
- *Inter-Task Mapping*: Directly map between source state and action spaces to target state and action space to derive a target policy.
- *New reinforcement learning methods*: Some algorithms inherently tackle transfer learning by design. Examples are available in [11].

Moreover, some approaches address specific subdomains such as multi-agent reinforcement learning [12], [13].

If source and target domain are not sufficiently similar, the performance gain through transfer learning may generally decrease. This problem is commonly known as negative transfer [6] and also occurs in the reinforcement learning context [11], [14]. Past research has commented on a suspected relation between negative transfer and overfitting [5].

Applications of transfer learning in reinforcement learning include training robots from simulated data [15], order dispatching in ridesharing platforms [16], deriving knowledge from sample-efficient 2D to cost-extensive 3D-environments [17] or gaming [18].

III. BACKGROUND

A. Deep Q-Networks

Motivated by the success of deep neural networks in computer vision at the time, [19] combined them with reinforcement learning techniques and managed to drastically outperform previous Q-learning approaches. Their so-called deep Q-network (DQN) utilizes a convolutional neural network (CNN) architecture with parameters θ that takes in raw pixel input s and outputs corresponding Q-values for the finite action space \mathcal{A} . For training, the model minimizes the mean squared error (MSE) between the current estimate $Q_\theta(s, a)$ and the target $r + \gamma \max_{a'} Q_\theta(s', a')$ as in

$$\nabla_\theta \text{MSE}(Q_\theta(s, a), r + \gamma \max_{a'} Q_\theta(s', a')) \quad (1)$$

As the resulting RL algorithm is off-policy the authors opted for using a replay memory, also known as experience replay buffer (ERB) [20], to store the last N generated experience samples. During training, batches are uniformly sampled from this buffer, thereby assigning the same importance to all experiences within it. This increases data efficiency as samples may be used more than once, and stabilizes training due to lower correlation between samples and an averaged behaviour distribution. However, the DQN implementation from [19] suffers from two problems: *Moving targets* and *overestimation bias*. The former arises as every update step that tries to improve the current Q-value estimate also changes the target. Consequently, convergence is difficult. The latter arises through the standard DQN setup where the same network is used for determining $\max_{a'}$ and calculating the target. In this case, the chosen Q-values will on average be larger than what they should be, leading to said overestimation bias [21]. To counter the problem of moving targets, [22] developed delayed DQNs. They introduce a second neural network θ' , called the *target network* to compute the targets. By not letting it receive gradient updates and only copying the parameters of the original network – the *model network* θ – at certain intervals, moving targets are alleviated.

B. Double Deep Q-Networks

To combat overestimation bias, [21] propose combining double Q learning [23] with DQNs naming the resulting method double DQN. By decoupling the action selection from its evaluation, they mitigate the risk of frequently picking overestimated Q-values as targets. For efficiency, the authors decided to use the target network to evaluate the action chosen by the model network. They support this decision, by arguing that the networks differ enough to make it unlikely that both overestimate the same values. This changes the target formulation from equation (1) to

$$\nabla_\theta \text{MSE}(Q_\theta(s, a), r + \gamma Q_{\theta'}(s', \arg\max_{a'} Q_\theta(s', a'))). \quad (2)$$

The authors were able to report significantly reduced overestimation and higher policy quality when comparing their approach to delayed DQNs on different Atari games [21].

C. Soft Q-learning

Reference [2] adapts Deep Q-Learning by prioritizing high entropy (i.e., unfamiliar) paths during training. The authors show that under minor preconditions, the fixed point iteration

$$Q_{soft}(s_t, a_t) \leftarrow r_t + \gamma \mathbb{E}_{s_{t+1} \sim p_s} [V_{soft}(s_{t+1})] \quad \forall s_t, a_t \quad (3)$$

$$V_{soft}(s_t) \leftarrow r_t + \alpha \int_{\mathcal{A}} \exp\left(\frac{Q_{soft}(s_t, \mathbf{a}')}{\alpha}\right) d\mathbf{a}' \quad \forall s_t \quad (4)$$

converge towards Q_{soft}^* and V_{soft}^* respectively to yield the optimal policy

$$\pi_{\text{MaxEnt}}^*(a_t|s_t) = \exp\left(\frac{Q_{soft}^*(s_t, a_t) - V_{soft}^*(s_t)}{\alpha}\right) \quad \forall s_t, a_t \quad (5)$$

for a heat parameter $\alpha > 0$ and an action space \mathcal{A} . The policy therefore resembles a softmax function and no longer acts greedily compared to Q-learning. Note that the distribution of \mathbf{a}' affects the (Lebesgue) integral in mapping (4) as \mathbf{a}' is a random variable. The authors also elaborate how their approach can handle continuous action spaces through approximation of the integral in mapping (4). However, as our use-case contains a low-dimensional, discrete action space, we could directly perform mapping (3) and (4). Moreover, in [24] the authors simplify mapping (4) to

$$V_{soft}(s_t) \leftarrow \mathbb{E}_{a_t \sim \pi} [Q_{soft}(s_t, a_t) - \log \pi(a_t|s_t)] \quad \forall s_t \quad (6)$$

As the entropy term $-\log \pi(a_t|s_t)$ is non-negative, mapping (6) implies that the policy prioritizes higher-entropy paths in that the policy is implicitly less confident. By scaling the entropy term by a fixed factor $\beta > 0$, we can further control how much the policy prioritizes high entropy paths. The term α is still important to control exploration. According to (5), for a given state s_t and actions a and a' , the likelihood to chose a' over a equals

$$\frac{\pi_{\text{MaxEnt}}^*(a'|s_t)}{\pi_{\text{MaxEnt}}^*(a|s_t)} = \exp\left(\frac{Q_{soft}^*(s_t, a') - Q_{soft}^*(s_t, a)}{\alpha}\right).$$

A larger choice of the heat parameter α , therefore, implies a more uniform policy π . (Hier auf Code aus Soft-Q Implementation verweisen -> In Training und Policy)

Intuitively, soft Q-Learning learns not only the single optimal, but multiple paths to achieve the goal due to its stochastic policy and targeted exploration strategy. As the Q-values in both Q-learning and soft Q-learning are subject to training noise, a Q-learning agent may often find itself on sub-optimal paths that it cannot follow. A soft Q-learning agent may be capable to recognize previously taken paths and still navigate itself to the goal.

Although research on soft Q-learning is scarce (searching "Soft Q-learning" on Google Scholar returned only 739 results), researchers have successfully applied it to real-world use-cases. For example, the Japanese billion-dollar company DeNA Co., Ltd. has applied soft Q-learning to distributed fleet control, reducing passenger waiting time by over 16.4% in

a realistic simulation architecture [25]. Other applications of soft Q-learning encompass pedestrian simulations [26] and text generation [27].

D. Prioritized Experience Replay Buffer

Prioritized replay [28] enhances reinforcement learning methods that utilize experience replay buffers. It aims to more frequently select important transitions from which the model can learn the most, instead of uniformly sampling from all experiences as is common for ERBs (see [19], [22]). This is done by treating the ERB as a priority queue that additionally tracks the current TD-error δ per experience. Training samples are then drawn via Thompson sampling. Since the TD-error for new interactions with the environment is still unknown, the authors suggest assigning a maximum priority to ensure that each observation affects training at least once. For efficiency, they further recommend only updating priorities for transitions after drawing them from the buffer and using them for training. To address the problems of reduced diversity in sample selection and overfitting that result from focusing on a small subset of high-priority experiences, [28] introduce a scaling parameter α to interpolate between uniform sampling ($\alpha = 0$) and greedy prioritization ($\alpha = 1$). The probability $P(i)$ of sampling experience i follows

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

with $p_i = |\delta_i| + \epsilon$ to always ensure a minimum base priority. Since the goal of training is to minimize the expected value of the TD-error, sampling by priorities introduces a bias that is corrected for through importance sampling. For each sampled experience, we obtain the corresponding importance sampling weight w_i according to

$$w_i = \left(\frac{1}{N} \frac{1}{P(i)}\right)^\beta$$

with N and β being the size of the buffer and the parameter for controlling the amount of correction, respectively. The authors report that gradually increasing β during training to 1 and setting α to 0.6 produces good results.

With their prioritized experience replay buffer (PERB), [28] increased performance over ERBs in 41 out of 49 of the Atari games used in [22]. Their technique has been widely adopted by various deep reinforcement learning algorithms, such as the infamous Rainbow [29] or the work of [30], and has successfully contributed to game-related reinforcement learning as in [31].

IV. METHODS

A. Environment

We measured and compared transfer learning capabilities on the popular NES game "Super Mario Bros.". The game consists of 32 levels in which the player has to control Super Mario through a parkour of obstacles (to reach a flag at the end). The player can chose from 256 distinct actions, all

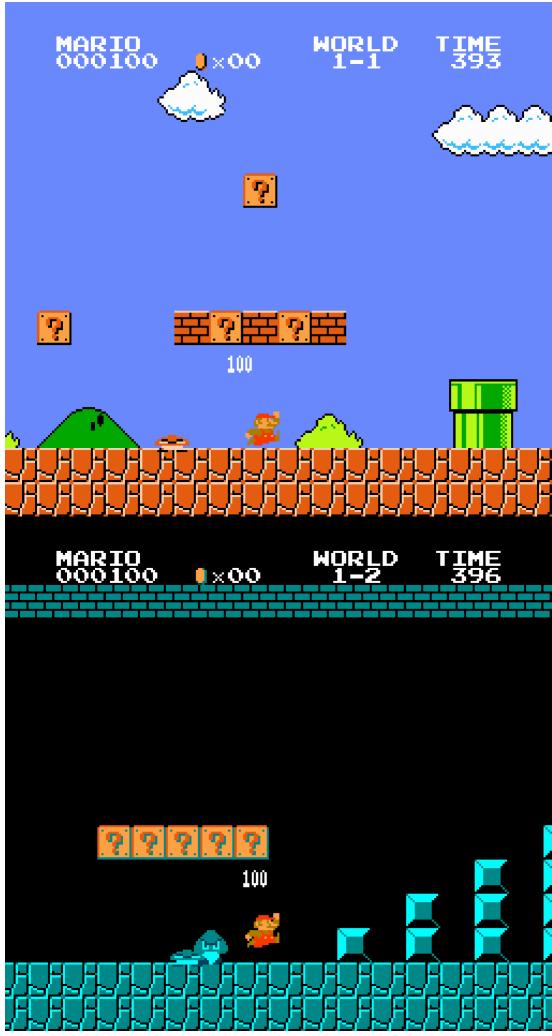


Fig. 1. Exemplary scenes from Super Mario Bros. levels 1-1 (top) and 1-2 (bottom).

related to moving right, left, jumping, throwing fire balls (only possible once Mario picked up a fire flower) or idling.

We relied on a ready-to-use implementation of Super Mario Bros.¹ that, among other ready-to-use reinforcement learning environments of various NES games², has served for previous reinforcement learning research [32], [33]. We perceive Super Mario Bros. as particularly well suited for reinforcement learning studies as its implementation is computationally cheap and it is well known and therefore comprehensible to other researchers. It is also very convenient to work with. We build the environment by calling

```
from nes_py.wrappers import JoypadSpace
import gym_super_mario_bros
from gym_super_mario_bros.actions
import SIMPLE_MOVEMENT
env = gym_super_mario_bros.make('SuperMarioBros-v0')
env = JoypadSpace(env, SIMPLE_MOVEMENT)
```

¹<https://pypi.org/project/gym-super-mario-bros/>

²<https://pypi.org/project/nes-py/>

and perform a (random) action via

```
action = env.action_space.sample()
state, reward, done, info = env.step(action)
env.close()
```

to receive a status, a reward, a confirmation whether the state is terminal and additional information about the game's state. Moreover, we render the game's current state by calling

```
env.render()
```

B. Reinforcement learning methods

To answer the first research question, we implemented soft Q-learning and chose DDQN as a comparative deterministic method. While DQN would be a more "vanilla" comparative method, we found its performance unsuitable for the problem at hand. DDQN extends DQN just slightly and performs sufficiently well. The pseudo code in Alg. 1 and Alg. 2 summarize our SQN and DDQN implementations. For numeric stability, we subtract the largest logit from all the logits so that the exponents in the softmax are always smaller or equal one. This is needed to prevent overflows:

```
logits = (q_values)/self.heat_param

# Subtract max for numeric stability
logits = logits - np.max(logits)

return logits
```

In [Appendix B], we display our implementations of the SQN's main characteristics – target calculation and action sampling.

C. Transfer learning approach

The transfer learning task consisted of learning to finish one level and leveraging the gained knowledge to finish another, which is characteristic for a starting-point-method. In reference to the taxonomy by [7], only the state spaces differ between levels and we only transfer the policy. Due to the exploratory nature of our research, we do not have a specific purpose for employing transfer learning.

As the various levels comprise similar visual assets and distinct challenges, they accurately resemble real world transfer learning use cases. We chose level 1-1 as the source domain and level 1-2 as the target domain. Because the two are among the most visually similar levels and the obstacles of level 1-1 also occur in level 1-2, we presumed that most information from learning the source domain could actually be transferred. Fig. 1 illustrates the visual similarities and differences between both levels.

For each DDQN and SQN, we trained two agents. The first agent learned on source domain until the cumulative reward per episode converged. Subsequently, we trained the agent further on the target domain. As a baseline, the other agent learned only on the target domain. We tracked the agents' win rates and cumulative rewards per episode. To pursue our second research question, we trained versions of the first agent on the target domain that had less training time on the source

Algorithm 1 DDQN Training Loop

```

1: Input: num_epochs, num_steps, decay_factor, batch_size,
   learning rate  $\eta$ , polyak parameter  $\tau$ , discount factor  $\gamma$ ,
   policy parameter  $\epsilon$ , PERB size  $N$  and parameters  $\alpha$  &  $\beta$ 
2: Initialize PERB with size  $N$  and fill it to its capacity
3: Initialize model network  $\theta$ 
4: Initialize delayed target network  $\theta'$ 
5:  $\theta' \leftarrow \theta$  ▷ Copy weights
6: for num_epochs do
7:    $\theta' \leftarrow (1 - \tau)\theta' + \tau\theta$  ▷ Polyak Averaging
8:   for num_steps do
9:     Interact with the environment ▷  $\epsilon$ -greedy action
10:    Store  $(s, a, r, s')$  in PERB with max prio
11:    if terminal then
12:      Reset environment
13:       $\epsilon \leftarrow \epsilon \cdot \text{decay\_factor}$  ▷ Anneal  $\epsilon$ 
14:    end if
15:  end for
16:  Sample batch_size transitions from the PERB (con-
   trolled by  $\alpha$ ) and obtain importance sampling weights
    $w$  (controlled by  $\beta$ ) ▷ see PERB implementation
17:   $y \leftarrow \begin{cases} r, & \text{if } s' = \text{terminal} \\ r + \gamma Q_{\theta'}(s', \arg\max_{a'} [Q_{\theta}(s', a')]), & \text{else} \end{cases}$ 
18:   $\delta \leftarrow \text{Huber}(y, Q_{\theta}(s, a))$ 
19:  Use  $\delta$  to update priors in PERB
20:  Increase  $\beta$  towards 1
21:   $\theta \leftarrow \theta + \eta[\nabla_{\theta}[w\delta]]$  ▷ SGD step with Adam
22: end for

```

domain. We essentially achieved this by loading the final weights from training on the first level:

```
model_network.load_weights("pathToWeights")
```

Before we compared SQN and DDQN with regard to transfer learning, we validated our SQN and DDQN implementations on the first level.

D. Preprocessing

The agents received the game state as a rescaled 84x84 grey-scale picture and drew from a restricted action space of five actions: (1) *idle*, (2) *move right*, (3) *jump right*, (4) *move right and throw a fire ball*, (5) *jump right and throw a fireball*. Because consecutive frames are highly correlated, we accelerated training by repeating each action over four frames and passing the corresponding states as a stacked 4x84x84 image. This way, we effectively skipped three out of four frames and informed the agent about moving objects. We also returned the cumulative reward over each four stacked frames after normalizing it into the range of $[-\frac{60}{600}, \frac{60}{600}]$. The implementation returns a positive reward equal to distance (in pixels) covered to the right, a penalty equal to time passed (in seconds) and a penalty of -15 for dying, but clips the total reward per frame into $[-15, 15]$. Due to choosing the discount factor as $\gamma = 0.9$, the absolute total discounted cumulative reward is therefore always less or equal $\frac{4 \cdot 15}{600} \cdot \frac{1}{1 - \gamma} = \frac{60}{600} \cdot 10 = 1$, which

Algorithm 2 SoftQ Training Loop

```

1: Input: num_epochs, num_steps, batch_size, heat_param,
   entropy_factor, learning rate  $\eta$ , polyak parameter  $\tau$ , dis-
   count factor  $\gamma$ , PERB size  $N$  and parameters  $\alpha$  &  $\beta$ 
2: Initialize PERB with size  $N$  and fill it to its capacity
3: Initialize model network  $\theta$ 
4: Initialize delayed target network  $\theta'$ 
5:  $\theta' \leftarrow \theta$  ▷ Copy weights
6: for num_epochs do
7:    $\theta' \leftarrow (1 - \tau)\theta' + \tau\theta$  ▷ Polyak Averaging
8:   for num_steps do
9:     Interact with the environment ▷ Entropy based
10:    Store  $(s, a, r, s')$  in PERB with max prio
11:    if terminal then
12:      Reset environment
13:    end if
14:  end for
15:  Sample batch_size transitions from the PERB (con-
   trolled by  $\alpha$ ) and obtain importance sampling weights
    $w$  (controlled by  $\beta$ ) ▷ see PERB implementation
16:   $\pi \leftarrow \text{softmax}(Q_{\theta'}(s') / \text{heat\_param})$ 
17:   $\text{entropy} \leftarrow \text{entropy\_factor} \cdot \sum [\pi \log \pi]$ 
18:   $V_{\text{soft}} \leftarrow \sum [\pi Q_{\theta'}(s')] - \text{entropy}$ 
19:   $y \leftarrow \begin{cases} r, & \text{if } s' = \text{terminal} \\ r + \gamma V_{\text{soft}}, & \text{else} \end{cases}$ 
20:   $\delta \leftarrow \text{Huber}(y, Q_{\theta}(s, a))$ 
21:  Use  $\delta$  to update priors in PERB
22:  Increase  $\beta$  towards 1
23:   $\theta \leftarrow \theta + \eta[\nabla_{\theta}[w\delta]]$  ▷ SGD step with Adam
24: end for

```

aims at stabilizing training. However, the expected discounted cumulative reward is around 0.6 and other reward scalings may be more suitable. We defined a wrapper skipping and stacking images and normalizing the rewards, and a separate function for resizing and greyscaling the image to exploit tensorflow's enhanced efficiency [Appendix A].

E. Hyperparameters

Starting from other practitioner's successful hyperparametrizations³, we manually and iteratively adapted network training hyperparameters.

Each epoch consisted of three environment steps and a single gradient step. To accelerate training, the model drew training samples through prioritized experience replay. For sampling observations, we used a sum tree approach [Appendix C]. Initially, we filled the buffer to its full size of 32,000 samples on a uniform policy (warm start). We further used the standard parameters $\alpha = 0.6$ and $\beta = 0.4$, and increased β towards one over the course of training.

Fig. 2 displays our network architecture. All convolutional layers used a kernel size of 3 and a stride of 1 without padding.

³<https://www.analyticsvidhya.com/blog/2021/06/playing-super-mario-bros-with-deep-reinforcement-learning/>

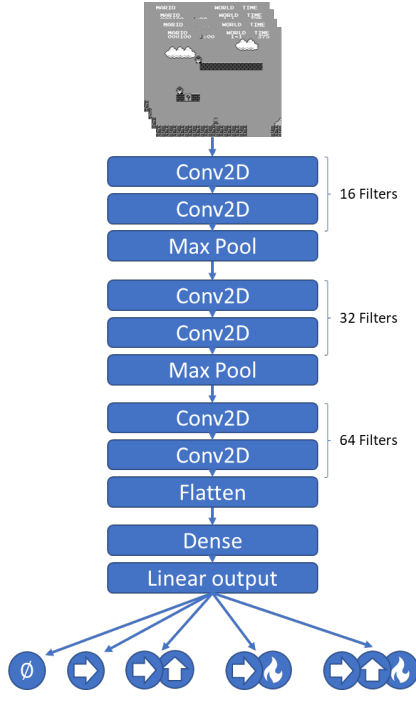


Fig. 2. Employed model network architecture

For pooling, we used a size of 2 and the dense layer contained 64 neurons. All models employed a learning rate of 0.00025 and a batch size of 64. As our PERB implementation's runtime scaled linearly in the batchsize and our submission deadline was firm, a larger batch size was infeasible.

Tab. I lists the respective hyperparameter choice for the DDQN and SQN. Both algorithms employed the same reward discount factor γ and Polyak averaging factor τ to prevent the problem of moving targets [22]. The initial exploration rate ϵ_{init} , the corresponding decay rate and the minimal exploration rate ϵ_{min} only relate to the DDQN. The heat parameter α and entropy scaling factor β only affect soft Q-learning. After varying both parameters in size, we chose them proportional to the reward scaling factor.

TABLE I
DDQN AND SQN HYPERPARAMETER VALUES

Hyperparameters	Algorithm	
	DDQN	SQN
γ	0.9	0.9
τ	$\frac{1}{3750}$	$\frac{1}{3750}$
ϵ_{init}	0.3	-
ϵ decay rate	0.999	-
ϵ_{min}	0.1	-
α	-	$\frac{1}{600}$
β	-	$\frac{1}{600}$

Fig. 3. SQN (top) and DDQN (bottom) running average over 100 episodes of cumulative reward on level 1-1. Values are multiplied by 600 to match the original reward distribution.

Fig. 4. SQN (top) and DDQN (bottom) running average over 100 episodes of cumulative wins on level 1-1.

F. Hardware configuration

We implemented all models in Python (v3.7.10 to v3.10.0) using the TensorFlow framework (v2.3.0 to v2.10.0) with GPU support. For training, we used two local PCs running Windows 10 (11) with a Ryzen 9 5900X 12x 3.700GHz (Intel Core i9-10900K 10x3,7Ghz) (CPU), 8GB RTX 3070 Ti FE (GPU) and 32GB DDR4-3200 (RAM).

V. RESULTS

To validate our SQN and DDQN implementations, we plotted their cumulative reward curves from training on the first level. Figs. 3 and 4 shows that both models successfully finished the level. The SQN performed better than the DDQN with a maximum win rate almost X times higher. The SQN also reached a higher average cumulative reward.

Figure 5 displays performance on level 1-2. Neither the SQN nor the DDQN seemed to benefit from transfer learning at all. Instead, the untrained versions seem to learn even faster than the pre-trained ones. Also, those versions of the policy with less training time do not appear to have learned any faster than the fully trained ones.

Goal:

- Summarize, what the data show
- simple relationships
- big picture trends
- cite figures / tables that present supporting data
- avoid simply repeating the numbers that are already available in figures / tables

Tips:

- Break into subsections with headings (if needed)
- Complement the information from the tables and figures
- give precise values that are not available in the figures
- report the percentage change of absolute values from the figures / tables
- repeat / highlight only the most important numbers - also talk about negative and control results
- reserve info about what you did for the methods section - in particular, why did something
- reserve comments of the meaning of the results for the discussion section

Contents:

- Num Episodes and Epochs
- Reward plots
- Nur auf section 2 - Dort jeweils ein Plot für DDQN und SQN, der sowohl pretrained als auch untrained enthält - auf section 1 reicht eigentlich ne erwähnung des finalen rewards

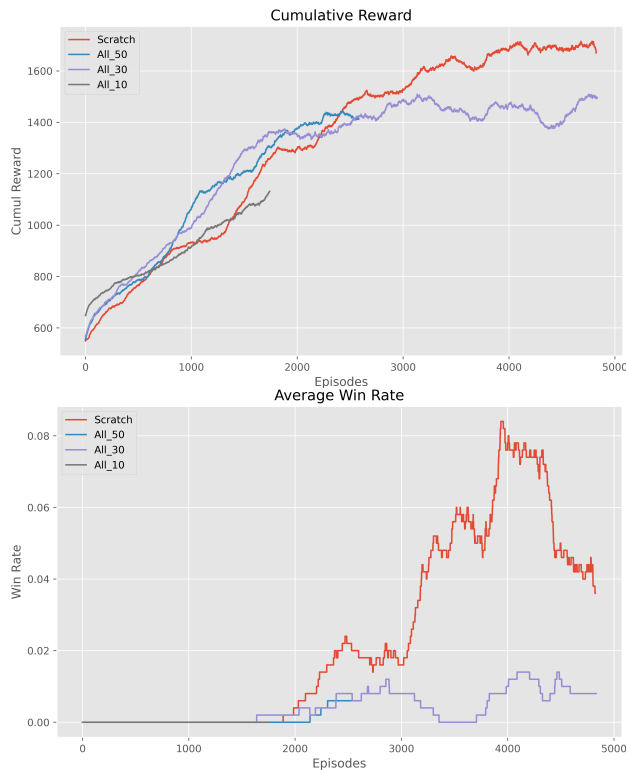


Fig. 5. SQN running averages over 500 episodes of cumulative reward (top) and wins (bottom) on level 1-2. The former are multiplied by 600 to match the original reward distribution.

im Text

-Win Rates

- Siehe reward plots

-Vergleich dieser Observations zwischen mit vs. ohne Transfer Learning, ohne vs. ohne, mit vs. mit. → Entweder innerhalb des gleichen Modells oder zwischen beiden verschiedenen

VI. DISCUSSION

Discussion here

As the cumulative reward of the untrained policy was the highest, neither the SQN nor the DDQN seemed to benefit from transfer learning at all. Hence, we could not observe any advantage or disadvantage of Soft Q-learning over a traditional approach for transfer learning. Moreover, we could also not observe any effect of training time on the source domain on performance on the target domain.

However, as we could not observe any (positive) effect through transfer learning, we assume that key mechanisms crucial for answering our research questions were not present. For example, the source and target domains may not have been compatible for classic transfer learning or we needed to apply a more sophisticated transfer learning approach. Just as well, we could have been in bad luck and required a larger number of observations to receive significant and reliable results. Therefore, we consider our research questions unanswered and call for future research to take them on.

Apart from other source and target domains, future research may also consider additional entropy-based and deterministic approaches to better generalize their findings. We hope our work can inspire other researchers investigate how to approach transfer learning both regarding methods and training time as answering either of our research questions may spare practitioners expensive resources as time and energy.

Ablauf:

- Answer questions asked
- Support your conclusion
- give possible mechanisms - strengths of our research -
- Defend your conclusion
- limitations of our research - Whats next - what needs to be confirmed - point out unanswered questions and future directions - Give the "big picture" take-home-message
- (human) implications
- tell the reader why they should care
- strong conclusion
- restate you main findings
- give a final take-home-message

Tips:

- Focus on what your data prove, not what you hoped it would prove
- focus on the limitations that matter, not generic limitations

VII. CONCLUSION

Conclusion here

ACKNOWLEDGMENTS

This research was funded by spilled lentils, Rohrbrüche, surging energy prices and two cursed Google Colab Pro subscriptions cancelling our trainings. We would like to thank our supervisor Leon for his thorough support (seriously), and Professor Bruni and Pipa for their wonderful informative lectures.

REFERENCES

- [1] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey, "Maximum Entropy Inverse Reinforcement Learning," 2008, p. 6. [Online]. Available: https://www.aaai.org/Papers/AAAI/2008/AAAI08-227.pdf?source=post_page
- [2] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, "Reinforcement Learning with Deep Energy-Based Policies," Jul. 2017. [Online]. Available: <http://arxiv.org/abs/1702.08165>
- [3] T. Haarnoja, V. Pong, A. Zhou, M. Dalal, P. Abbeel, and S. Levine, "Composable Deep Reinforcement Learning for Robotic Manipulation," 2018, arXiv:1803.06773 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1803.06773>
- [4] M. Nemecek and R. Parr, "Policy Caches with Successor Features," in *Proceedings of the 38th International Conference on Machine Learning*. PMLR, 2021, pp. 8025–8033. [Online]. Available: <https://proceedings.mlr.press/v139/nemecek21a.html>
- [5] S. Gamrian and Y. Goldberg, "Transfer Learning for Related Reinforcement Learning Tasks via Image-to-Image Translation," in *Proceedings of the 36th International Conference on Machine Learning*. PMLR, 2019, pp. 2063–2072, iSSN: 2640-3498. [Online]. Available: <https://proceedings.mlr.press/v97/gamrian19a.html>

- [6] Z. Wang, Z. Dai, B. Póczos, and J. Carbonell, "Characterizing and Avoiding Negative Transfer," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, 2019, pp. 11 285–11 294. [Online]. Available: <https://doi.org/10.1109/CVPR.2019.01155>
- [7] A. Lazaric, "Transfer in Reinforcement Learning: A Framework and a Survey," in *Reinforcement Learning*, M. Wiering and M. van Otterlo, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, vol. 12, pp. 143–173. [Online]. Available: http://link.springer.com/10.1007/978-3-642-27645-3_5
- [8] M. E. Taylor and P. Stone, "Transfer Learning for Reinforcement Learning Domains: A Survey," *The Journal of Machine Learning Research*, vol. 10, pp. 1633–1685, 2009.
- [9] Y. Zhan and M. E. Taylor, "Online Transfer Learning in Reinforcement Learning Domain," 2015, p. 8. [Online]. Available: <https://www.aaai.org/ocs/index.php/FSS/FSS15/paper/view/11646>
- [10] Z. Zhu, K. Lin, A. K. Jain, and J. Zhou, "Transfer Learning in Deep Reinforcement Learning: A Survey," 2022, arXiv:2009.07888 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/2009.07888>
- [11] L. Torrey and J. Shavlik, "Transfer learning," in *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*. IGI global, 2010, pp. 242–264. [Online]. Available: <https://doi.org/10.4018/978-1-60566-766-9.ch011>
- [12] G. Boutsioukis, I. Partalas, and I. Vlahavas, "Transfer Learning in Multi-Agent Reinforcement Learning Domains," in *Recent Advances in Reinforcement Learning*, ser. Lecture Notes in Computer Science, S. Sanner and M. Hutter, Eds. Berlin, Heidelberg: Springer, 2012, pp. 249–260. [Online]. Available: https://doi.org/10.1007/978-3-642-29946-9_25
- [13] F. L. D. Silva and A. H. R. Costa, "A Survey on Transfer Learning for Multiagent Reinforcement Learning Systems," *Journal of Artificial Intelligence Research*, vol. 64, pp. 645–703, 2019. [Online]. Available: <https://doi.org/10.1613/jair.1.11396>
- [14] S. A. H. Minoofam, A. Bastanfard, and M. R. Keyvanpour, "TRCLA: A Transfer Learning Approach to Reduce Negative Transfer for Cellular Learning Automata," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–10, 2021. [Online]. Available: <https://doi.org/10.1109/TNNLS.2021.3106705>
- [15] W. Zhao, J. P. Queralta, and T. Westerlund, "Sim-to-Real Transfer in Deep Reinforcement Learning for Robotics: a Survey," in *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. Canberra, ACT, Australia: IEEE, 2020, pp. 737–744. [Online]. Available: <https://doi.org/10.1109/SSCI47803.2020.9308468>
- [16] R. Wan, S. Zhang, C. Shi, S. Luo, and R. Song, "Pattern Transfer Learning for Reinforcement Learning in Order Dispatching," 2021, arXiv:2105.13218 [cs]. [Online]. Available: <http://arxiv.org/abs/2105.13218>
- [17] L. A. Celiberto, J. P. Matsuura, R. L. de Mantaras, and R. A. Bianchi, "Using Transfer Learning to Speed-Up Reinforcement Learning: A Cased-Based Approach," in *2010 Latin American Robotics Symposium and Intelligent Robotics Meeting*. Sao Bernardo do Campo: IEEE, 2010, pp. 55–60. [Online]. Available: <https://doi.org/10.1109/LARS.2010.24>
- [18] K. Shao, Y. Zhu, and D. Zhao, "StarCraft Micromanagement With Reinforcement Learning and Curriculum Transfer Learning," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 3, no. 1, pp. 73–84, 2019. [Online]. Available: <https://doi.org/10.1109/TETCI.2018.2823329>
- [19] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing Atari with Deep Reinforcement Learning," 2013. [Online]. Available: <http://arxiv.org/abs/1312.5602>
- [20] L.-J. Lin, "Self-improving reactive agents based on reinforcement learning, planning and teaching," *Machine Learning*, vol. 8, no. 3, pp. 293–321, 1992. [Online]. Available: <https://doi.org/10.1007/BF00992699>
- [21] H. van Hasselt, A. Guez, and D. Silver, "Deep Reinforcement Learning with Double Q-Learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016. [Online]. Available: <https://doi.org/10.1609/aaai.v30i1.10295>
- [22] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015. [Online]. Available: <https://doi.org/10.1038/nature14236>
- [23] H. van Hasselt, "Double Q-learning," in *Advances in Neural Information Processing Systems*, vol. 23. Curran Associates, Inc., 2010. [Online]. Available: <https://papers.nips.cc/paper/2010/hash/091d584fcd301b442654dd8c23b3fc9-Abstract.html>
- [24] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor," 2018, arXiv:1801.01290 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1801.01290>
- [25] T. Oda and Y. Tachibana, "Distributed Fleet Control with Maximum Entropy Deep Reinforcement Learning," in *Machine Learning for Intelligent Transportation Systems, Poster Submission*, 2018, p. 7.
- [26] F. Martínez-Gil, M. Lozano, I. García-Fernández, P. Romero, D. Serra, and R. Sebastián, "Using Inverse Reinforcement Learning with Real Trajectories to Get More Trustworthy Pedestrian Simulations," *Mathematics*, vol. 8, no. 9, p. 1479, 2020. [Online]. Available: <https://doi.org/10.3390/math8091479>
- [27] H. Guo, B. Tan, Z. Liu, E. P. Xing, and Z. Hu, "Text Generation with Efficient (Soft) Q-Learning," 2021. [Online]. Available: <http://arxiv.org/abs/2106.07704>
- [28] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized Experience Replay," in *4th International Conference on Learning Representations, ICLR 2016, May 2-4, 2016, Conference Track Proceedings*, San Juan, Puerto Rico, 2016. [Online]. Available: <http://arxiv.org/abs/1511.05952>
- [29] M. Hessel, J. Modayil, H. v. Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver, "Rainbow: Combining Improvements in Deep Reinforcement Learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018. [Online]. Available: <https://doi.org/10.1609/aaai.v32i1.11796>
- [30] S. Kaptrowski, G. Ostrovski, J. Quan, R. Munos, and W. Dabney, "Recurrent Experience Replay in Distributed Reinforcement Learning," in *7th International Conference on Learning Representations, ICLR 2019, May 6-9, 2019*, New Orleans, LA, USA, 2019.
- [31] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver, "Mastering Atari, Go, chess and shogi by planning with a learned model," *Nature*, vol. 588, no. 7839, pp. 604–609, 2020. [Online]. Available: <https://doi.org/10.1038/s41586-020-03051-4>
- [32] D. Chen, Y. Wang, and W. Gao, "Combining a gradient-based method and an evolution strategy for multi-objective reinforcement learning," *Applied Intelligence*, vol. 50, no. 10, pp. 3301–3317, 2020. [Online]. Available: <https://doi.org/10.1007/s10489-020-01702-7>
- [33] D. G. LeBlanc and G. Lee, "General Deep Reinforcement Learning in NES Games," *Proceedings of the Canadian Conference on Artificial Intelligence*, 2021. [Online]. Available: <https://doi.org/10.21428/594757db.8472938b>

APPENDIX

A. Preprocessing

1) Frame Skipping:

```
class SkipObs(gym.Wrapper):
    """Wrapper class that repeats the chosen action for a given amount of frames
    """

    def __init__(self, env=None, skip: int=4, reward_scale_factor: int=1):
        """Initialize the wrapper

        Arguments:
            env (gym): Environment the wrapper should be applied to
            skip (int): Number of frames to skip
            reward_scale_factor (int): Factor to scale the rewards with
        """

        super(SkipObs, self).__init__(env)

        self.skip = skip
        self.reward_scale_factor = reward_scale_factor

    def step(self, action: int):
        """Function to overwrite the environments step function

        Arguments:
            action (int): The chosen action

        Returns:
            obs (nd.array): Last observation from the environment
            total_reward (float): Total reward from all steps (including skipped ones)
            done (boolean): Whether the episode is finished or not
            info (dict): Dictionary containing information about the environments state
        """

        total_reward = 0.0
        done = None

        # Apply the same action to "skip" many frames and calculate the total reward
        for _ in range(self.skip):
            obs, reward, done, info = self.env.step(action)

            # Reward scaling
            total_reward += reward / self.reward_scale_factor

            if done:
                break

        return obs, total_reward, done, info
```

2) Wrapper Combination:

```
# For environment creation
def make_env(level: str, movement_type: list, num_skip: int=4, num_stack: int=4, reward_scale_factor: int=1
):
    """Function to apply all wrappers to the environment

    Arguments:
        level (str): Level to create
        movement_type (list): Defines the joypad space e.g. allow movement to the right only
        num_skip (int): Number of frames to skip
        num_stack (int): Number of frames stack
        reward_scale_factor (int): Factor to scale the rewards with

    Returns:
        env (gym): Wrapped environment
    """

    # Create Env from given level
    env = gym_super_mario_bros.make(level)
```

```

# Select joypad space
env = JoypadSpace(env=env, actions=movement_type)

# Skip "num_skip" frames
env = SkipObs(env=env, skip=num_skip, reward_scale_factor=reward_scale_factor)

# Wrapper that stacks "num_stack" frames resulting in observations being of size (num_stack,240,256,3)
env = FrameStack(env=env, num_stack=num_stack)

return env

```

3) Image Resizing and Greyscaling:

```

@tf.function
def preprocess_obs(obs, resize_env):
    """ Function to preprocess the observations before they are fed into the network
    1. Rescale the pixel values between [0,1]
    2. Convert from RGB to greyscale
    3. Resize to given resize_env value
    4. Transpose stacked images such that consecutive ones are within the channel dimension

    Arguments:
        obs (gym.wrappers.LazyFrames): Stacked observations from the environment
        resize_env (tuple): Tuple containing the new size to which the observation should be adjusted

    Returns:
        preprocessed_obs (tf.tensor): Preprocessed observation
    """

    preprocessed_obs = tf.convert_to_tensor(obs/255)
    preprocessed_obs = tf.image.rgb_to_grayscale(preprocessed_obs)
    preprocessed_obs = tf.image.resize(preprocessed_obs, size=resize_env)
    preprocessed_obs = tf.transpose(tf.squeeze(preprocessed_obs), perm=[1,2,0])

    return preprocessed_obs

```

B. Soft Q-Learning

1) Target Calculation:

```

for s, a, r, s_prime, terminal in dataset:
    # Get q-values from target network
    q_prime = target_network(s_prime)

    logits = q_prime/policy.heat_param
    logits = logits - tf.math.reduce_max(logits, axis=1, keepdims=True)
    pi = tf.math.softmax(logits, axis=1)

    if np.sum(np.isnan(pi)) > 0:
        pi = tf.where(tf.equal(tf.reduce_max(q_prime, axis=1, keep_dims=True), q_prime),
                      tf.constant(1, shape=pi.shape),
                      tf.constant(0, shape=pi.shape))

    # Calculate the entropy of the policy pi in state s
    entropy = entropy_factor * tf.reduce_sum(pi*tf.math.log(pi + 0.001), axis=1)

    # Calculate v soft' according to Equation (3) in Haarnoja et al. 2018
    # https://arxiv.org/abs/1801.01290
    V_soft_prime = tf.reduce_sum(pi*(q_prime), axis=1) - entropy

```

2) Policy:

```

# Policy Class
# Doesn't use any tf functions as numpy turned out to be faster, even compared to graph mode
class Policy():
    """Policy used to sample an action given the q-values of the current state
    """

    def __init__(self, heat_param: float=1):
        """Initialize the policy

        Arguments:
            heat_param (float): Heat parameter that controls the exploration of the policy

```

```

"""

self.heat_param = heat_param

def __call__(self, q_values):
    """Sample action according to the q-values

    Arguments:
        q_values (tf.Tensor): q-values of the current state given by the model network

    Returns:
        action (int): To-be-performed action
    """

    p = self.calc_probability(q_values)
    action = np.random.choice(len(q_values), p=p)

    return action

def calc_probability(self, q_values):
    """Calculate the entropy-based probabilities for the given q-values

    Arguments:
        q_values (tf.Tensor): q-values of the current state given by the model network

    Returns:
        p (nd.array): Array of probabilities
    """

    # Calculate logits
    logits = self.calc_logits(q_values)

    p = np.exp(logits)
    p /= np.sum(p)

    # If p contains nan (occurs if one value is much larger than the others)
    # then replace the max with 1 and all others with zero
    # --> deterministic sampling for numeric stability
    if np.sum(np.isnan(p)) > 0:
        p = np.where(q_values == np.max(q_values), 1, 0)

    return p

def calc_logits(self, q_values):
    """Calculate logits for the given q-values scaled by the heat parameter

    Arguments:
        q_values (tf.Tensor): q-values of the current state given by the model network

    Returns:
        logits (nd.array): Logits of the q values
    """

    logits = (q_values)/self.heat_param

    # Subtract max for numeric stability
    logits = logits - np.max(logits)

    return logits

```

C. Prioritized Experience Replay Buffer

1) Weight Adjustment:

```

def adjust_priority(self, priority: float):
    """Function to adjust a priority

    Arguments:
        priority (float): To-be-adjusted priority

    Returns:

```

```

        adjusted_priority (float): Adjusted priority
    """

    # Add min prio to current prio and raise sum to the power of alpha
    adjusted_priority = np.power(priority + self.min_priority, self.alpha)

    return adjusted_priority

```

2) Sampling:

```

def sample(self, num_samples: int):
    """Function for sampling

    Arguments:
        num_samples (int): Number of samples to draw from the buffer

    Returns:
        dataset (tf.data.Dataset): Dataset containing the sampled experiences
        sampled_idx (list): List of the sampled experience indices
        importance_sampling_weights (list): List of the corresponding importance sampling weights
    """

    sampled_idx = []
    importance_sampling_weights = []
    sample_no = 0

    # Sample "num_samples" indices from the sum tree
    while sample_no < num_samples:

        # Uniformly sample between 0 and the summed priority (stored in the root node)
        # used for retrieving an experience sample index (i.e., sampling according to priorities)
        sample_val = np.random.uniform(0, self.base_node.value)
        sampled_node = retrieve(sample_val, self.base_node)

        if sampled_node.idx < self.available_samples - 1:

            sampled_idx.append(sampled_node.idx)

            # Divide the sampled nodes' priority by the summed priority and calculate the importance
            # sampling weight
            p = sampled_node.value / self.base_node.value
            importance_sampling_weights.append((self.available_samples + 1) * p)

            sample_no += 1

    # Apply the beta factor and normalize so that the maximum importance sampling weight < 1
    importance_sampling_weights = np.array(importance_sampling_weights)
    importance_sampling_weights = np.power(importance_sampling_weights, -self.beta)
    importance_sampling_weights = importance_sampling_weights / np.max(importance_sampling_weights)

    obs, actions, rewards, obs_prime, terminal = [], [], [], [], []

    for idx in sampled_idx:
        # Append the sampled experience components from the buffer to the respective lists
        obs.append(self.buffer[idx][self.obs_idx])
        actions.append(self.buffer[idx][self.action_idx])
        rewards.append(self.buffer[idx][self.reward_idx])
        obs_prime.append(self.buffer[idx][self.obs_prime_idx])
        terminal.append(self.buffer[idx][self.terminal_idx])

    dataset = tf.data.Dataset.from_tensor_slices((obs, actions, rewards, obs_prime, terminal)).batch(
        num_samples)

    return dataset, sampled_idx, importance_sampling_weights

```