# Sample-Efficient Deep Reinforcement Learning via Episodic Backward Update

**Su Young Lee,    Sungik Choi,    Sae-Young Chung**
School of Electrical Engineering, KAIST, Republic of Korea
`{suyoung.l, si_choi, schung}@kaist.ac.kr`

## Abstract

We propose Episodic Backward Update (EBU) – a novel deep reinforcement learning algorithm with a direct value propagation. In contrast to the conventional use of the experience replay with uniform random sampling, our agent samples a whole episode and successively propagates the value of a state to its previous states. Our computationally efficient recursive algorithm allows sparse and delayed rewards to propagate directly through all transitions of the sampled episode. We theoretically prove the convergence of the EBU method and experimentally demonstrate its performance in both deterministic and stochastic environments. Especially in 49 games of Atari 2600 domain, EBU achieves the same mean and median human normalized performance of DQN by using only 5% and 10% of samples, respectively.

## 1   Introduction

Deep reinforcement learning (DRL) has been successful in many complex environments such as the Arcade Learning Environment [2] and Go [18]. Despite DRL's impressive achievements, it is still impractical in terms of sample efficiency. To achieve human-level performance in the Arcade Learning Environment, Deep $Q$-Network (DQN) [14] requires 200 million frames of experience for training which corresponds to 39 days of gameplay in real-time. Clearly, there is still a tremendous gap between the learning process of humans and that of deep reinforcement learning agents. This problem is even more crucial for tasks such as autonomous driving, where we cannot risk many trials and errors due to the high cost of samples.

One of the reasons why DQN suffers from such low sample efficiency is the sampling method from the replay memory. In many practical problems, an RL agent observes sparse and delayed rewards. There are two main problems when we sample one-step transitions uniformly at random. **(1)** We have a low chance of sampling a transition with a reward for its sparsity. The transitions with rewards should always be updated to assign credits for actions that maximize the expected return. **(2)** In the early stages of training when all values are initialized to zero, there is no point in updating values of one-step transitions with zero rewards if the values of future transitions with nonzero rewards have not been updated yet. Without the future reward signals propagated, the sampled transition will always be trained to return a zero value.

In this work, we propose Episodic Backward Update (EBU) to present solutions for the problems raised above. When we observe an event, we scan through our memory and seek for the past event that caused the later one. Such an episodic control method is how humans normally recognize the cause and effect relationship [10]. Inspired by this, we can solve the first problem **(1)** by sampling transitions in an episodic manner. Then, we can be assured that at least one transition with a non-zero reward is used for the value update. We can solve the second problem **(2)** by updating the values of transitions in a backward manner in which the transitions were made. Afterward, we can perform an

efficient reward propagation without any meaningless updates. This method faithfully follows the principle of dynamic programming.

As mentioned by the authors of DQN, updating correlated samples in a sequence is vulnerable to overestimation. In Section 3, we deal with this issue by adopting a diffusion factor to mediate between the learned values from the future transitions and the current sample reward. In Section 4, we theoretically prove the convergence of our method for both deterministic and stochastic MDPs. In Section 5, we empirically show the superiority of our method on 2D MNIST Maze Environment and the 49 games of Atari 2600 domain. Especially in 49 games of the Atari 2600 domain, our method requires only 10M frames to achieve the same mean human-normalized score reported in Nature DQN [14], and 20M frames to achieve the same median human-normalized score. Remarkably, EBU achieves such improvements with a comparable amount of computation complexity by only modifying the target generation procedure for the value update from the original DQN.

## 2 Background

The goal of reinforcement learning (RL) is to learn the optimal policy that maximizes the expected sum of rewards in the environment that is often modeled as a Markov decision process (MDP) $M = (\mathcal{S}, \mathcal{A}, P, R)$. $\mathcal{S}$ denotes the state space, $\mathcal{A}$ denotes the action space, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ denotes the transition probability distribution, and $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ denotes the reward function. $Q$-learning [22] is one of the most widely used methods to solve RL tasks. The objective of $Q$-learning is to estimate the state-action value function $Q(s, a)$, or the $Q$-function, which is characterized by the Bellman optimality equation. $Q^*(s_t, a) = \mathbb{E}[r_t + \gamma \max_{a'} Q^*(s_{t+1}, a')]$.

There are two major inefficiencies of the traditional on-line $Q$-learning. First, each experience is used only once to update the $Q$-function. Secondly, learning from experiences in a chronologically forward order is much more inefficient than learning in a chronologically backward order, because the value of $s_{t+1}$ is required to update the value of $s_t$. Experience replay [12] is proposed to overcome these inefficiencies. After observing a transition $(s_t, a_t, r_t, s_{t+1})$, the agent stores the transition into its replay buffer. In order to learn the $Q$-values, the agent samples transitions from the replay buffer.

In practice, the state space $\mathcal{S}$ is extremely large, therefore it is impractical to tabularize the $Q$-values of all state-action pairs. Deep $Q$-Network [14] overcomes this issue by using deep neural networks to approximate the $Q$-function. DQN adopts experience replay to use each transition for multiple updates. Since DQN uses a function approximator, consecutive states output similar $Q$-values. If DQN updates transitions in a chronologically backward order, often overestimation errors cumulate and degrade the performance. Therefore, DQN does not sample transitions in a backward order, but uniformly at random. This process breaks down the correlations between consecutive transitions and reduces the variance of updates.

There have been a variety of methods proposed to improve the performance of DQN in terms of stability, sample efficiency, and runtime. Some methods propose new network architectures. The dueling network architecture [21] contains two streams of separate $Q$-networks to estimate the value functions and the advantage functions. Neural episodic control [16] and model-free episodic control [5] use episodic memory modules to estimate the state-action values. RUDDER [1] introduces an LSTM network with contribution analysis for an efficient return decomposition. Ephemeral Value Adjustments (EVA) [7] combines the values of two separate networks, where one is the standard DQN and another is a trajectory-based value network.

Some methods tackle the uniform random sampling replay strategy of DQN. Prioritized experience replay [17] assigns non-uniform probability to sample transitions, where greater probability is assigned for transitions with higher temporal difference (TD) error. Inspired by Lin's backward use of replay memory, some methods try to aggregate TD values with Monte-Carlo (MC) returns. $Q(\lambda)$ [23], $Q^*(\lambda)$ [6] and Retrace($\lambda$) [15] modify the target values to allow the on-policy samples to be used interchangeably for on-policy and off-policy learning. Count-based exploration method combined with intrinsic motivation [3] takes a mixture of one-step return and MC return to set up the target value. Optimality Tightening [8] applies constraints on the target using the values of several neighboring transitions. Simply by adding a few penalty terms to the loss, it efficiently propagates reliable values to achieve fast convergence.
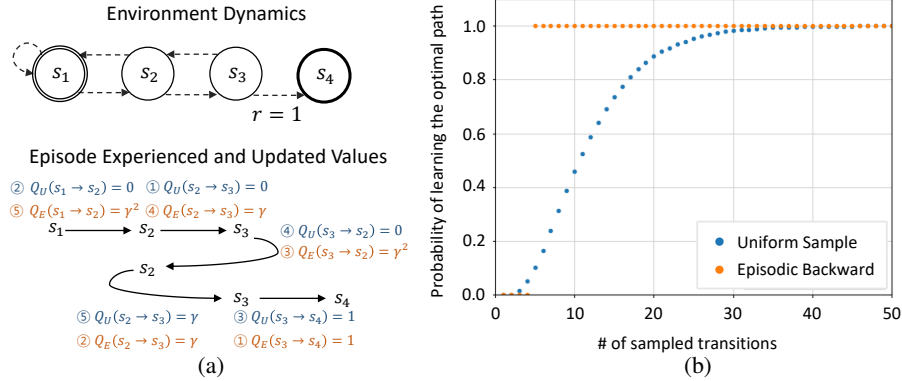
Figure 1: A motivating example where uniform sampling method fails but EBU does not. (a): A simple navigation domain with 4 states and a single rewarded transition. Circled numbers indicate the order of sample updates. $Q_U$ and $Q_E$ stand for the $Q$-values learned by the uniform random sampling method and the EBU method respectively. (b): The probability of learning the optimal path $(s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4)$ after updating the $Q$-values with sample transitions.

Our goal is to improve the sample efficiency of deep reinforcement learning by making a simple yet effective modification. Without a single change of the network structure, training schemes, and hyperparameters of the original DQN, we only modify the target generation method. Instead of using a limited number of transitions, our method samples a whole episode from the replay memory and propagates the values sequentially throughout the entire transitions of the sampled episode in a backward manner. By using a temporary backward $Q$-table with a diffusion coefficient, our novel algorithm effectively reduces the errors generated from the consecutive updates of correlated states.

## 3 Proposed Methods

### 3.1 Episodic Backward Update for Tabular $Q$-Learning

Let us imagine a simple tabular MDP with a single rewarded transition (Figure 1, (a)), where an agent can only take one of the two actions: *'left'* and *'right'*. In this example, $s_1$ is the initial state, and $s_4$ is the terminal state. A reward of 1 is gained only when the agent reaches the terminal state and a reward of 0 is gained from any other transitions. To make it simple, assume that we have only one episode stored in the experience memory: $(s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4)$. The $Q$-values of all transitions are initialized to zero. With a discount $\gamma \in (0, 1)$, the optimal policy is to take the action *'right'* in all states. When sampling transitions uniformly at random as Nature DQN, the key transitions $(s_1 \rightarrow s_2)$, $(s_2 \rightarrow s_3)$ and $(s_3 \rightarrow s_4)$ may not be sampled for updates. Even when those transitions are sampled, there is no guarantee that the update of the transition $(s_3 \rightarrow s_4)$ is done before the update of $(s_2 \rightarrow s_3)$. We can speed up the reward propagation by updating all transitions within the episode in a backward manner. Such a recursive update is also computationally efficient.

We can calculate the probability of learning the optimal path $(s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow s_4)$ as a function of the number of sample transitions trained. With the tabular Episodic Backward Update stated in Algorithm 1, which is a special case of Lin's algorithm [11] with recency parameter $\lambda = 0$, the agent can figure out the optimal policy just after 5 updates of $Q$-values. However, we see that the uniform sampling method requires more than 40 transitions to learn the optimal path with probability close to 1 (Figure 1, (b)).

Note that this method differs from the standard $n$-step $Q$-learning [22]. In $n$-step $Q$-learning, the number of future steps for the target generation is fixed as $n$. However, our method considers $T$ future values, where $T$ is the length of the sampled episode. $N$-step $Q$-learning takes a max operator at the $n$-th step only, whereas our method takes a max operator at every iterative backward step which can propagate high values faster. To avoid exponential decay of the $Q$-value, we set the learning rate $\alpha = 1$ within the single episode update.

3

---

**Algorithm 1** Episodic Backward Update for Tabular $Q$-Learning (single episode, tabular)

---

1: Initialize the $Q$- table $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ with all-zero matrix.
   $Q(s, a) = 0$ for all state action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$.
2: Experience an episode $E = \{(s_1, a_1, r_1, s_2), \dots, (s_T, a_T, r_T, s_{T+1})\}$
3: **for** $t = T$ to $1$ **do**
4: $\quad Q(s_t, a_t) \leftarrow r_t + \gamma \max_{a'} Q(s_{t+1}, a')$
5: **end for**

---

---

**Algorithm 2** Episodic Backward Update

---

1: **Initialize**: replay memory $D$ to capacity $N$, on-line action-value function $Q(\cdot; \boldsymbol{\theta})$, target action-value function $\hat{Q}(\cdot; \boldsymbol{\theta}^-)$
2: **for** episode = 1 to $M$ **do**
3: $\quad$ **for** $t = 1$ to Terminal **do**
4: $\quad\quad$ With probability $\epsilon$ select a random action $a_t$, otherwise select $a_t = \operatorname{argmax}_a Q(s_t, a; \boldsymbol{\theta})$
5: $\quad\quad$ Execute action $a_t$, observe reward $r_t$ and next state $s_{t+1}$
6: $\quad\quad$ Store transition $(s_t, a_t, r_t, s_{t+1})$ in $D$
7: $\quad\quad$ Sample a random episode $E = \{\boldsymbol{S}, \boldsymbol{A}, \boldsymbol{R}, \boldsymbol{S'}\}$ from $D$, set $T = \text{length}(E)$
8: $\quad\quad$ Generate a temporary target $Q$-table, $\tilde{Q} = \hat{Q}(\boldsymbol{S'}, \cdot; \boldsymbol{\theta}^-)$
9: $\quad\quad$ Initialize the target vector $\boldsymbol{y} = \text{zeros}(T)$, $\boldsymbol{y}_T \leftarrow \boldsymbol{R}_T$
10: $\quad\quad$ **for** $k = T - 1$ to $1$ **do**
11: $\quad\quad\quad$ $\tilde{Q}[\boldsymbol{A}_{k+1}, k] \leftarrow \beta \boldsymbol{y}_{k+1} + (1 - \beta)\tilde{Q}[\boldsymbol{A}_{k+1}, k]$
12: $\quad\quad\quad$ $\boldsymbol{y}_k \leftarrow \boldsymbol{R}_k + \gamma \max_a \tilde{Q}[a, k]$
13: $\quad\quad$ **end for**
14: $\quad\quad$ Perform a gradient descent step on $(\boldsymbol{y} - Q(\boldsymbol{S}, \boldsymbol{A}; \boldsymbol{\theta}))^2$ with respect to $\boldsymbol{\theta}$
15: $\quad\quad$ Every $C$ steps reset $\hat{Q} = Q$
16: $\quad$ **end for**
17: **end for**

---

There are some other multi-step methods that converge to the optimal state-action value function, such as $Q(\lambda)$ and $Q^*(\lambda)$. However, our algorithm neither cuts trace of trajectories as $Q(\lambda)$, nor requires the parameter $\lambda$ to be small enough to guarantee convergence as $Q^*(\lambda)$. We present a detailed discussion on the relationship between EBU and other multi-step methods in Appendix F.

## 3.2 Episodic Backward Update for Deep $Q$-Learning[1]

Directly applying the backward update algorithm to deep reinforcement learning is known to show highly unstable results due to the high correlation of consecutive samples. We show that the fundamental ideas of the tabular version of the backward update algorithm may be applied to its deep version with just a few modifications. The full algorithm introduced in Algorithm 2 closely resembles that of Nature DQN [14]. Our contributions lie in the recursive backward target generation with a diffusion factor $\beta$ (starting from line number 7 of Algorithm 2), which prevents the overestimation errors from correlated states cumulating.

Instead of sampling transitions uniformly at random, we make use of all transitions within the sampled episode $E = \{\boldsymbol{S}, \boldsymbol{A}, \boldsymbol{R}, \boldsymbol{S'}\}$. Let the sampled episode start with a state $S_1$, and contain T transitions. Then $E$ can be denoted as a set of four length-$T$ vectors: $\boldsymbol{S} = \{S_1, S_2, \dots, S_T\}$; $\boldsymbol{A} = \{A_1, A_2, \dots, A_T\}$; $\boldsymbol{R} = \{R_1, R_2, \dots, R_T\}$ and $\boldsymbol{S'} = \{S_2, S_3, \dots, S_{T+1}\}$. The temporary target $Q$-table, $\tilde{Q}$ is an $|\mathcal{A}| \times T$ matrix which stores the target $Q$-values of all states $\boldsymbol{S'}$ for all valid actions, where $\mathcal{A}$ is the action space of the MDP. Therefore, the $j$-th column of $\tilde{Q}$ is a column vector that contains $\hat{Q}(S_{j+1}, a; \boldsymbol{\theta}^-)$ for all valid actions $a \in \mathcal{A}$, where $\hat{Q}$ is the target $Q$-function parametrized by $\boldsymbol{\theta}^-$.

After the initialization of the temporary $Q$-table, we perform a recursive backward update. Adopting the backward update idea, one element $\tilde{Q}[\boldsymbol{A}_{k+1}, k]$ in the $k$-th column of the $\tilde{Q}$ is replaced using the next transition's target $\boldsymbol{y}_{k+1}$. Then $\boldsymbol{y}_k$ is estimated as the maximum value of the newly modified $k$-th column of $\tilde{Q}$. Repeating this procedure in a recursive manner until the start of the episode, we can

---

[1]The code is available at `https://github.com/suyoung-lee/Episodic-Backward-Update`

successfully apply the backward update algorithm for a deep $Q$-network. The process is described in detail with a supplementary diagram in Appendix E.

We are using a function approximator, and updating correlated states in a sequence. As a result, we observe overestimated values propagating and compounding through the recursive $\max$ operations. We solve this problem by introducing the diffusion factor $\beta$. By setting $\beta \in (0, 1)$, we can take a weighted sum of the new backpropagated value and the pre-existing value estimate. One can regard $\beta$ as a learning rate for the temporary $Q$-table, or as a level of *'backwardness'* of the update. This process stabilizes the learning process by exponentially decreasing the overestimation error. Note that Algorithm 2 with $\beta = 1$ is identical to the tabular backward algorithm stated in Algorithm 1. When $\beta = 0$, the algorithm is identical to episodic one-step DQN. The role of $\beta$ is investigated in detail with experiments in Section 5.3.

### 3.3 Adaptive Episodic Backward Update for Deep $Q$-Learning

The optimal diffusion factor $\beta$ varies depending on the type of the environment and the degree of how much the network is trained. We may further improve EBU by developing an adaptive tuning scheme for $\beta$. Without increasing the sample complexity, we propose an adaptive, single actor and multiple learner version of EBU. We generate $K$ learner networks with different diffusion factors, and a single actor to output a policy. For each episode, the single actor selects one of the learner networks in a regular sequence. Each learner is trained in parallel, using the same episode sampled from a shared experience replay. Even with the same training data, all learners show different interpretations of the sample based on the different levels of trust in backwardly propagated values. We record the episode scores of each learner during training. After every fixed step, we synchronize all the learner networks with the parameters of a learner network with the best training score. This adaptive version of EBU is presented as a pseudo-code in Appendix A. In Section 5.2, we compare the two versions of EBU, one with a constant $\beta$ and another with an adaptive $\beta$.

## 4 Theoretical Convergence

### 4.1 Deterministic MDPs

We prove that Episodic Backward Update with $\beta \in (0, 1)$ defines a contraction operator, and converges to the optimal $Q$-function in finite and deterministic MDPs.

**Theorem 1.** *Given a finite, deterministic and tabular MDP $M = (\mathcal{S}, \mathcal{A}, P, R)$, the Episodic Backward Update algorithm in Algorithm 3 converges to the optimal Q-function w.p. 1 as long as*

- *The step size satisfies the Robbins-Monro condition;*

- *The sample trajectories are finite in lengths $l$: $\mathbb{E}[l] < \infty$;*

- *Every (state, action) pair is visited infinitely often.*

We state the proof of Theorem 1 in Appendix G. Furthermore, even in stochastic environments, we can guarantee the convergence of the episodic backward algorithm for a sufficiently small $\beta$.

### 4.2 Stochastic MDPs

**Theorem 2.** *Given a finite, tabular and stochastic MDP $M = (\mathcal{S}, \mathcal{A}, P, R)$, define $R_{\max}^{\mathrm{sto}}(s, a)$ as the maximal return of trajectory that starts from state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$. In a similar way, define $r_{\min}^{\mathrm{sto}}(s, a)$ and $r_{\mathrm{mean}}^{\mathrm{sto}}(s, a)$ as the minimum and mean of possible reward by selecting action $a$ in state $s$. Define $\mathcal{A}_{\mathrm{sub}}(s) = \{a' \in \mathcal{A} | Q^*(s, a') < \max_{a \in \mathcal{A}} Q^*(s, a)\}$ as the set of suboptimal actions in state $s \in \mathcal{S}$. Define $\mathcal{A}_{\mathrm{opt}}(s) = \mathcal{A} \backslash \mathcal{A}_{\mathrm{sub}}(s)$. Then, under the conditions of Theorem 1, and*

$$\beta \leq \inf_{s \in \mathcal{S}} \inf_{a' \in \mathcal{A}_{\mathrm{sub}}(s)} \inf_{a \in \mathcal{A}_{\mathrm{opt}}(s)} \frac{Q^*(s, a) - Q^*(s, a')}{R_{\max}^{\mathrm{sto}}(s, a') - Q^*(s, a')}, \tag{1}$$

$$\beta \leq \inf_{s \in \mathcal{S}} \inf_{a' \in \mathcal{A}_{\mathrm{sub}}(s)} \inf_{a \in \mathcal{A}_{\mathrm{opt}}(s)} \frac{Q^*(s, a) - Q^*(s, a')}{r_{\mathrm{mean}}^{\mathrm{sto}}(s, a) - r_{\min}^{\mathrm{sto}}(s, a)}, \tag{2}$$

*the Episodic Backward Update algorithm in Algorithm 3 converges to the optimal Q-function w.p. 1.*
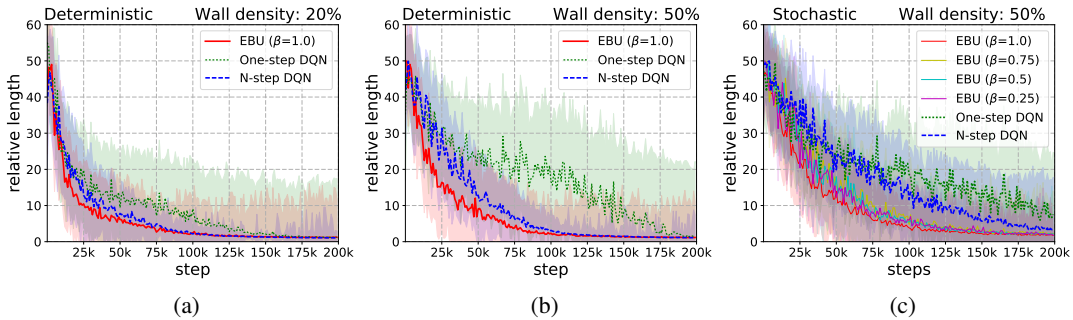
Figure 2: (a) & (b): Median of 50 relative lengths of EBU and baselines. EBU outperforms other baselines significantly in the low sample regime and for high wall density. (c): Median relative lengths of EBU and other baseline algorithms in MNIST maze with stochastic transitions.

The main intuition of this theorem is that $\beta$ acts as a learning rate of the backward target therefore mitigates the collision between the max operator and stochastic transitions.

## 5 Experimental Results

### 5.1 2D MNIST Maze (Deterministic/Stochastic MDPs)

We test our algorithm in the 2D Maze Environment. Starting from the initial position at $(0, 0)$, the agent has to navigate through the maze to reach the goal position at $(9, 9)$. To minimize the correlation between neighboring states, we use the MNIST dataset [9] for the state representation. The agent



Figure 3: 2D MNIST Maze

receives the coordinates of the position in two MNIST images as the state representation. The training environments are 10 by 10 mazes with randomly placed walls. We assign a reward of 1000 for reaching the goal, and a reward of -1 for bumping into a wall. A wall density indicates the probability of having a wall at each position. For each wall density, we generate 50 random mazes with different wall locations. We train a total of 50 independent agents, one for each maze over 200,000 steps. The performance metric, relative length is defined as $l_{\mathrm{rel}} = l_{\mathrm{agent}}/l_{\mathrm{oracle}}$, which is the ratio between the length of the agent's path $l_{\mathrm{agent}}$ and the length of the ground truth shortest path $l_{\mathrm{oracle}}$ to reach the goal. The details of the hyperparameters and the network structure are described in Appendix D.

We compare EBU to uniform random sampling one-step DQN and $n$-step DQN. For $n$-step DQN, we set the value of $n$ as the length of the episode. Since all three algorithms eventually achieve median relative lengths of 1 at the end of the training, we report the relative lengths at 100,000 steps in Table 1. One-step DQN performs the worst in all configurations, implying the inefficiency of uniform sampling update in environments with sparse and delayed rewards. As the wall density increases, it becomes more important for the agent to learn the correct decisions at bottleneck positions. $N$-step DQN shows the best performance with a low wall density, but as the wall density increases, EBU significantly outperforms $n$-step DQN.

In addition, we run experiments with stochastic transitions. We assign 10% probability for each side action for all four valid actions. For example, when an agent takes an action '*up*', there is a 10% chance of transiting to the left state, and 10% chance of transiting to the right state. In Figure 9 (c), we see that the EBU agent outperforms the baselines in the stochastic environment as well.
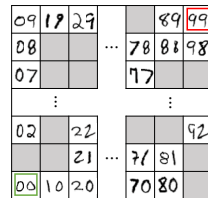
Table 1: Relative lengths (Mean & Median) of 50 deterministic MNIST Maze after 100,000 steps

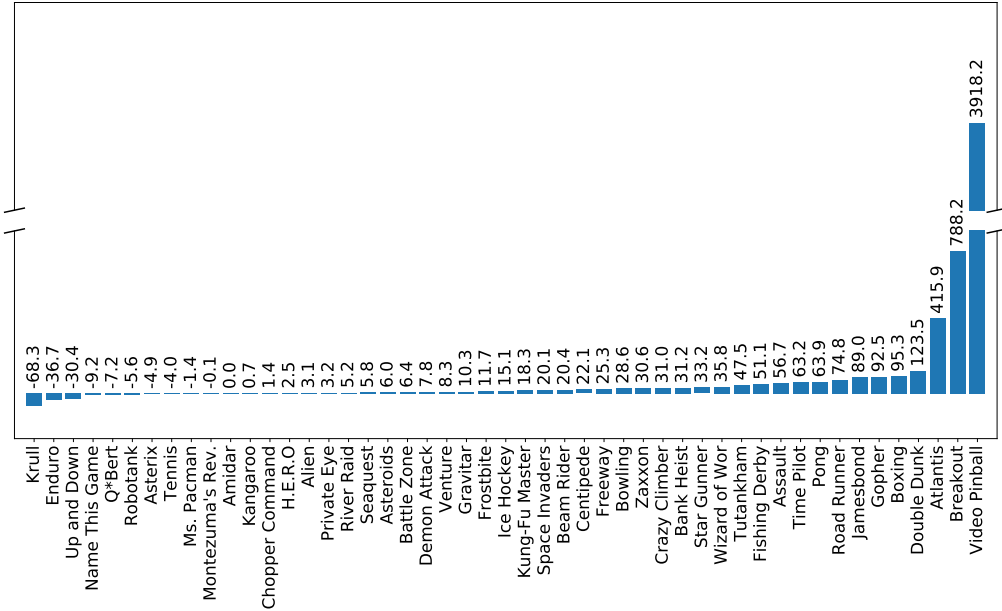| Wall density | EBU ($\beta = 1.0$) | | One-step DQN | | $N$-step DQN | |
|---|---|---|---|---|---|---|
| 20% | 5.44 | 2.42 | 14.40 | 9.25 | **3.26** | **2.24** |
| 30% | **8.14** | **3.03** | 25.63 | 21.03 | 8.88 | 3.32 |
| 40% | **8.61** | **2.52** | 25.45 | 22.71 | 8.96 | 3.50 |
| 50% | **5.51** | **2.34** | 22.36 | 16.62 | 11.32 | 3.12 |

6

Figure 4: Relative score of adaptive EBU (4 random seeds) compared to Nature DQN (8 random seeds) in percents (%) both trained for 10M frames.

## 5.2 49 Games of Atari 2600 Environment (Deterministic MDPs)

The Arcade Learning Environment [2] is one of the most popular RL benchmarks for its diverse set of challenging tasks. We use the same set of 49 Atari 2600 games, which was evaluated in Nature DQN paper [14].

We select $\beta = 0.5$ for EBU with a constant diffusion factor. For adaptive EBU, we train $K = 11$ parallel learners with diffusion factors $0.0, 0.1, \ldots$, and $1.0$. We synchronize the learners at the end of each epoch (0.25M frames). We compare our algorithm to four baselines: Nature DQN [14], Prioritized Experience Replay (PER) [17], Retrace($\lambda$) [15] and Optimality Tightening (OT) [8]. We train EBU and baselines for 10M frames (additional 20M frames for adaptive EBU) on 49 Atari games with the same network structure, hyperparameters, and evaluation methods used in Nature DQN. The choice of such a small number of training steps is made to investigate the sample efficiency of each algorithm following [16, 8]. We report the mean result from 4 random seeds for adaptive EBU and 8 random seeds for all other baselines. Detailed specifications for each baseline are described in Appendix D.

First, we show the improvement of adaptive EBU over Nature DQN at 10M frames for all 49 games in Figure 4. To compare the performance of an agent to its baseline's, we use the following relative score, $\frac{\text{Score}_{\text{Agent}} - \text{Score}_{\text{Baseline}}}{\max\{\text{Score}_{\text{Human}}, \text{Score}_{\text{Baseline}}\} - \text{Score}_{\text{Random}}}$ [21]. This measure shows how well an agent performs a task compared to the task's level of difficulty. EBU ($\beta = 0.5$) and adaptive EBU outperform Nature DQN in 33 and 39 games out of 49 games, respectively. The large amount of improvements in games such as "Atlantis," "Breakout," and "Video Pinball" highly surpass minor failings in few games.

We use human-normalized score, $\frac{\text{Score}_{\text{Agent}} - \text{Score}_{\text{Random}}}{|\text{Score}_{\text{Human}} - \text{Score}_{\text{Random}}|}$ [20], which is the most widely used metric to make an apple-to-apple comparison in the Atari domain. We report the mean and the median human-normalized scores of the 49 games in Table 2. The result signifies that our algorithm outperforms the baselines in both the mean and median of the human-normalized scores. PER and Retrace($\lambda$) do not show a lot of improvements for a small number of training steps as 10M frames. Since OT has to calculate the $Q$-values of neighboring states and compare them to generate the penalty term, it requires about 3 times more training time than Nature DQN. However, EBU performs iterative episodic updates using the temporary $Q$-table that is shared by all transitions in the episode, EBU has almost the same computational cost as that of Nature DQN.
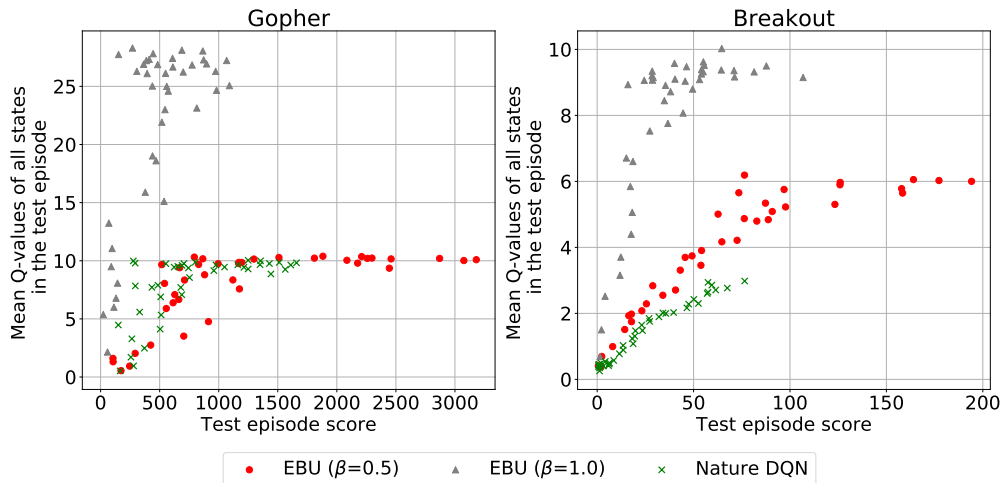
Figure 5: Episode scores and average $Q$-values of all state-action pairs in "Gopher" and "Breakout".

The most significant result is that EBU ($\beta = 0.5$) requires only 10M frames of training to achieve the mean human-normalized score reported in Nature DQN, which is trained for 200M frames. Although 10M frames are not enough to achieve the same median score, adaptive EBU trained for 20M frames achieves the median normalized score. These results signify the efficacy of backward value propagation in the early stages of training. Raw scores for all 49 games are summarized in Appendix B. Learning curves of adaptive EBU for all 49 games are reported in Appendix C.

Table 2: Summary of training time and human-normalized performance. Training time refers to the total time required to train 49 games of 10M frames using a single NVIDIA TITAN Xp for a single random seed. We use multi-GPUs to train learners of adaptive EBU in parallel. (*) The result of OT differs from the result reported in [8] due to different evaluation methods (i.e. not limiting the maximum number of steps for a test episode and taking maximum score from random seeds). (**) We report the scores of Nature DQN (200M) from [14].

| Algorithm (frames) | Training Time (hours) | Mean (%) | Median (%) |
|---|---|---|---|
| EBU ($\beta = 0.5$) (10M) | 152 | 253.55 | 51.55 |
| EBU (adaptive $\beta$) (10M) | 203 | 275.78 | 63.80 |
| Nature DQN (10M) | 138 | 133.95 | 40.42 |
| PER (10M) | 146 | 156.57 | 40.86 |
| Retrace($\lambda$) (10M) | 154 | 93.77 | 41.99 |
| OT (10M)* | 407 | 162.66 | 49.42 |
| EBU (adaptive $\beta$) (20M) | 450 | 347.99 | 92.50 |
| Nature DQN (200M)** | - | 241.06 | 93.52 |

## 5.3 Analysis on the Role of the Diffusion Factor $\beta$

In this section, we make comparisons between our own EBU algorithms. EBU ($\beta = 1.0$) works the best in the MNIST Maze environment because we use MNIST images for the state representation to allow consecutive states to exhibit little correlation. However, in the Atari domain, consecutive states are often different in a scale of few pixels only. As a consequence, EBU ($\beta = 1.0$) underperforms EBU ($\beta = 0.5$) in most of the Atari games. In order to analyze this phenomenon, we evaluate the $Q$-values learned at the end of each training epoch. We report the test episode score and the corresponding mean $Q$-values of all transitions within the test episode (Figure 5). We notice that the EBU ($\beta = 1.0$) is trained to output highly overestimated $Q$-values compared to its actual return. Since the EBU method performs recursive `max` operations, EBU outputs higher (possibly overestimated) $Q$-values than Nature DQN. This result indicates that sequentially updating correlated states with

overestimated values may destabilize the learning process. However, this result clearly implies that EBU ($\beta = 0.5$) is relatively free from the overestimation problem.

Next, we investigate the efficacy of using an adaptive diffusion factor. In Figure 6, we present how adaptive EBU adapts its diffusion factor during the course of training in "Breakout". In the early stage of training, the agent barely succeeds in breaking a single brick. With a high $\beta$ close to 1, values can be directly propagated from the rewarded state to the state where the agent has to bounce the ball up. Note that the performance of adaptive EBU follows that of EBU ($\beta = 1.0$) up to about 5M frames. As the training proceeds, the agent encounters more rewards and various trajectories that may cause overestimation. As a consequence, we discover that the agent anneals the diffusion factor to a lower value of 0.5. The trend of how the diffusion factor adapts differs from game to game. Refer to the diffusion factor curves for all 49 games in Appendix C to check how adaptive EBU selects the best diffusion factor.
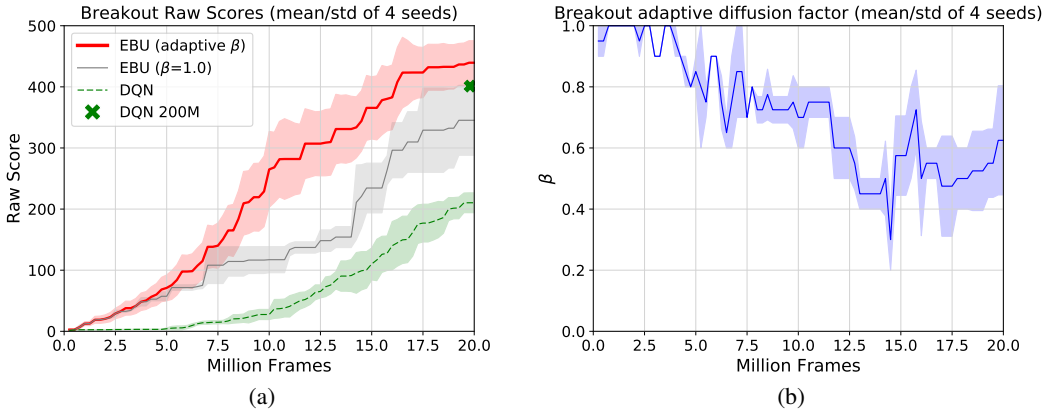


Figure 6: (a) Test scores in "Breakout". Mean and standard deviation from 4 random seeds are plotted. (b) Adaptive diffusion factor of adaptive EBU in "Breakout".

## 6 Conclusion

In this work, we propose Episodic Backward Update, which samples transitions episode by episode, and updates values recursively in a backward manner. Our algorithm achieves fast and stable learning due to its efficient value propagation. We theoretically prove the convergence of our method, and experimentally show that our algorithm outperforms other baselines in many complex domains, requiring only about 10% of samples. Since our work differs from DQN only in terms of the target generation, we hope that we can make further improvements by combining with other successful deep reinforcement learning methods.

## References

[1] Arjona-Medina, J. A., Gillhofer, M., Widrich, M., Unterthiner, T., and Hochreiter, S. RUDDER: Return decomposition for delayed rewards. arXiv preprint arXiv:1806.07857, 2018.

[2] Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. Journal of Artificial Intelligence Research, 47:253-279, 2013.

[3] Bellemare, M. G., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. Unifying count-based exploration and intrinsic motivation. In Advances in Neural Information Processing Systems (NIPS), 1471-1479, 2016.

[4] Bertsekas, D. P., and Tsitsiklis, J. N. Neuro-Dynamic Programming. Athena Scientific, 1996.

[5] Blundell, C., Uria, B., Pritzel, A., Li, Y., Ruderman, A., Leibo, J. Z, Rae, J.,Wierstra, D., and Hassabis, D. Modelfree episodic control. arXiv preprint arXiv:1606.04460, 2016.

[6] Harutyunyan, A., Bellemare, M. G., Stepleton, T., and Munos, R. Q($\lambda$) with off-policy corrections. In International Conference on Algorithmic Learning Theory (ALT), 305-320, 2016.

[7] Hansen, S., Pritzel, A., Sprechmann, P., Barreto, A., and Blundell, C. Fast deep reinforcement learning using online adjustments from the past. In Advances in Neural Information Processing Systems (NIPS), 10590–10600, 2018

[8] He, F. S., Liu, Y., Schwing, A. G., and Peng, J. Learning to play in a day: Faster deep reinforcement learning by optimality tightening. In International Conference on Learning Representations (ICLR), 2017.

[9] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. In the Institute of Electrical and Electronics Engineers (IEEE), 86, 2278-2324, 1998.

[10] Lengyel, M., and Dayan, P. Hippocampal Contributions to Control: The Third Way. In Advances in Neural Information Processing Systems (NIPS), 889-896, 2007.

[11] Lin, L-J. Programming Robots Using Reinforcement Learning and Teaching. In Association for the Advancement of Artificial Intelligence (AAAI), 781-786, 1991.

[12] Lin, L-J. Self-improving reactive agents based on reinforcement learning, planning and teaching. Machine Learning, 293-321, 1992.

[13] Melo, F. S. Convergence of Q-learning: A simple proof, Institute Of Systems and Robotics, Tech. Rep, 2001.

[14] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. Nature, 518(7540):529-533, 2015.

[15] Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. G. Safe and efficient off-policy reinforcement learning. In Advances in Neural Information Processing Systems (NIPS), 1046-1054, 2016.

[16] Pritzel, A., Uria, B., Srinivasan, S., Puig-'domenech, A., Vinyals, O., Hassabis, D., Wierstra, D., and Blundell, C. Neural Episodic Control. In International Conference on Machine Learning (ICML), 2827-2836, 2017.

[17] Schaul, T., Quan, J., Antonoglou, I., and Silver, D. Prioritized Experience Replay. In International Conference on Learning Representations (ICLR), 2016.

[18] Silver, D., Huang, A., Maddison C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I., Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., and Hassabis, D. Mastering the game of Go with deep neural networks and tree search. Nature, 529:484-489, 2016.

[19] Sutton, R. S., and Barto, A. G. Reinforcement Learning: An Introduction. MIT Press, 1998.

[20] van Hasselt, H., Guez, A., and Silver, D. Deep Reinforcement Learning with Double Q-learning. In Association for the Advancement of Artificial Intelligence (AAAI), 2094-2100, 2016.

[21] Wang, Z., Schaul, T., Hessel, M., van Hasselt, H., Lanctot, M., and de Freitas, N. Dueling Network Architectures for Deep Reinforcement Learning. In International Conference on Machine Learning (ICML), 1995-2003, 2016.

[22] Watkins., C. J. C. H. Learning from delayed rewards. Ph.D. thesis, University of Cambridge England, 1989.

[23] Watkins., C. J. C. H., and Dayan, P. Q-learning. Machine Learning, 272-292, 1992.

# Appendix A  Episodic Backward Update with an adaptive diffusion factor

---

**Algorithm 3** Adaptive Episodic Backward Update

---

1: **Initialize**: replay memory $D$ to capacity $N$, $K$ on-line action-value function $Q_1(\cdot;\boldsymbol{\theta_1}),\ldots,Q_K(\cdot;\boldsymbol{\theta_K})$, $K$ target action-value function $\hat{Q}_1(\cdot;\boldsymbol{\theta_1^-}),\ldots,\hat{Q}_K(\cdot;\boldsymbol{\theta_K^-})$, training score recorder $TS = \mathrm{zeros}(K)$, diffusion factors $\beta_1,\ldots,\beta_K$ for each learner network
2: **for** episode = 1 to $M$ **do**
3:     Select $Q_{\mathrm{actor}} = Q_i$ as the actor network for the current episode, where $i = (\text{episode} - 1)\%K + 1$
4:     **for** $t = 1$ to Terminal **do**
5:         With probability $\epsilon$ select a random action $a_t$
6:         Otherwise select $a_t = \mathrm{argmax}_a Q_{\mathrm{actor}}(s_t, a)$
7:         Execute action $a_t$, observe reward $r_t$ and next state $s_{t+1}$
8:         Store transition $(s_t, a_t, r_t, s_{t+1})$ in $D$
9:         Add training score for the current learner $TS[i] + = r_t$
10:         Sample a random episode $E = \{\boldsymbol{S}, \boldsymbol{A}, \boldsymbol{R}, \boldsymbol{S'}\}$ from $D$, set $T = \mathrm{length}(E)$
11:         **for** $j = 1$ to $K$ (this loop is processed in parallel) **do**
12:             Generate temporary target $Q$-table, $\tilde{Q}_j = \hat{Q}_i\left(\boldsymbol{S'}, \cdot; \boldsymbol{\theta_j^-}\right)$
13:             Initialize target vector $\boldsymbol{y} = \mathrm{zeros}(T)$, $\boldsymbol{y}_T \leftarrow \boldsymbol{R}_T$
14:             **for** $k = T - 1$ to 1 **do**
15:                 $\tilde{Q}_j\left[\boldsymbol{A}_{k+1}, k\right] \leftarrow \beta_j \boldsymbol{y}_{k+1} + (1 - \beta_j)\tilde{Q}_j\left[\boldsymbol{A}_{k+1}, k\right]$
16:                 $\boldsymbol{y}_k \leftarrow \boldsymbol{R}_k + \gamma \max_a \tilde{Q}_j[a, k]$
17:             **end for**
18:             Perform a gradient descent step on $\left(\boldsymbol{y} - Q_j\left(\boldsymbol{S}, \boldsymbol{A}; \boldsymbol{\theta_j}\right)\right)^2$ with respect to $\boldsymbol{\theta_j}$
19:         **end for**
20:         Every $C$ steps reset $\hat{Q}_1 = Q_1,\ldots,\hat{Q}_K = Q_K$
21:     **end for**
22:     Every $B$ steps synchronize all learners with the best training score, $b = \mathrm{argmax}_k TS[k]$. $Q_1(\cdot;\boldsymbol{\theta_1}) = Q_b(\cdot;\boldsymbol{\theta_b}),\ldots,Q_K(\cdot;\boldsymbol{\theta_K}) = Q_b(\cdot;\boldsymbol{\theta_b})$ and $\hat{Q}_1(\cdot;\boldsymbol{\theta_1^-}) = \hat{Q}_b(\cdot;\boldsymbol{\theta_b^-}),\ldots,\hat{Q}_K(\cdot;\boldsymbol{\theta_K^-}) = \hat{Q}_b(\cdot;\boldsymbol{\theta_b^-})$. Reset the training score recorder $TS = \mathrm{zeros}(K)$.
23: **end for**

---

# Appendix B  Raw scores of all 49 games.

Table 3: Raw scores after 10M frames of training. Mean scores from 4 random seeds are reported for adaptive EBU. 8 random seeds are used for all other baselines. We use the results at Nature DQN paper to report the scores at 200M frames. We run their code (`https://github.com/deepmind/dqn`) to report scores for 10M frames. Due to the use of different random seeds, the result of Nature DQN at 10M frames may be better than that of Nature DQN at 200M frames in some games. Bold texts indicate the best score out of the 5 results trained for 10M frames.

| Training Frames | 10M | | | | | | 20M | 200M |
|---|---|---|---|---|---|---|---|---|
| | EBU($\beta$=0.5) | Adap. EBU | DQN | PER | Retrace($\lambda$) | OT | Adap. EBU | Nature DQN |
| Alien | 708.08 | 894.15 | 690.32 | 1026.96 | 708.29 | **1078.67** | 1225.36 | 3069.00 |
| Amidar | 117.94 | 124.63 | 125.42 | 167.63 | 182.68 | **220.00** | 209.96 | 739.50 |
| Assault | **4109.18** | 3676.95 | 2426.94 | 2720.69 | 2989.05 | 2499.23 | 3943.23 | 3359.00 |
| Asterix | 1898.12 | 2533.27 | **2936.54** | 2218.54 | 1798.54 | 2592.50 | 3221.25 | 6012.00 |
| Asteroids | 1002.17 | **1402.43** | 654.99 | 993.50 | 886.92 | 985.88 | 2378.84 | 1629.00 |
| Atlantis | 61708.75 | 87944.38 | 20666.84 | 35663.83 | **98182.81** | 57520.00 | 141226.00 | 85641.00 |
| Bank heist | 359.62 | **459.42** | 234.70 | 312.96 | 223.50 | 407.42 | 680.43 | 429.70 |
| Battle zone | 20627.73 | 24748.50 | 22468.75 | 20835.74 | **30128.36** | 20400.48 | 30502.53 | 26300.00 |
| Beam rider | 5628.99 | 4785.27 | 3682.92 | 4586.07 | 4093.76 | **5889.54** | 6634.43 | 6846.00 |
| Bowling | 52.02 | **102.89** | 65.23 | 42.74 | 42.62 | 53.45 | 113.75 | 42.40 |
| Boxing | 55.95 | **72.69** | 37.28 | 4.64 | 6.76 | 60.89 | 96.35 | 71.80 |
| Breakout | 174.76 | **265.62** | 28.36 | 164.22 | 171.86 | 75.00 | 443.34 | 401.20 |
| Centipede | 4651.28 | **8389.16** | 6207.30 | 4385.41 | 5986.16 | 5277.79 | 8389.16 | 8309.00 |
| Chopper Command | 1196.67 | 1294.45 | 1168.67 | 1344.24 | 1353.76 | **1615.00** | 1909.23 | 6687.00 |
| Crazy Climber | 65329.63 | **94135.04** | 74410.74 | 53166.47 | 64598.21 | 92972.08 | 103780.15 | 114103.00 |
| Demon Attack | 7924.14 | **8368.16** | 7772.39 | 4446.03 | 6450.84 | 6872.04 | 9099.16 | 9711.00 |
| Double Dunk | -16.19 | **-14.12** | -17.94 | -15.62 | -15.81 | -15.92 | -12.78 | -18.10 |
| Enduro | 415.59 | 326.45 | 516.10 | 308.75 | 208.10 | **615.05** | 410.95 | 301.80 |
| Fishing Derby | -39.13 | **-15.85** | -65.53 | -78.49 | -75.74 | -69.66 | 9.22 | -0.80 |
| Freeway | 19.07 | **23.71** | 16.24 | 9.35 | 15.26 | 14.63 | 34.36 | 30.30 |
| Frostbite | 437.92 | 966.23 | 466.02 | 536.00 | 825.00 | **2452.75** | 1760.15 | 328.30 |
| Gopher | 3318.50 | **3634.67** | 1726.52 | 1833.67 | 3410.75 | 2869.08 | 5611.30 | 8520.00 |
| Gravitar | 294.58 | **450.18** | 193.55 | 319.79 | 272.08 | 263.54 | 611.79 | 306.70 |
| H.E.R.O. | 3089.90 | 3398.55 | 2767.97 | 3052.04 | 3079.43 | **10698.25** | 4308.23 | 19950.00 |
| Ice Hockey | -4.71 | **-2.96** | -4.79 | -7.73 | -6.13 | -5.79 | -2.96 | -1.60 |
| Jamesbond | 391.67 | **519.52** | 183.35 | 421.46 | 436.25 | 325.21 | 1043.66 | 576.70 |
| Kangaroo | 535.83 | 731.13 | 709.88 | **782.50** | 538.33 | 708.33 | 2018.83 | 6740.00 |
| Krull | 7587.24 | 8733.52 | **24109.14** | 6642.58 | 6346.40 | 7468.70 | 10016.72 | 3805.00 |
| Kung-Fu Master | 20578.33 | **26069.68** | 21951.72 | 18212.89 | 18815.83 | 22211.25 | 30387.78 | 23270.00 |
| Montezuma's Revenge | 0.00 | 0.00 | **3.95** | 0.43 | 0.00 | 0.00 | 0.00 | 0.00 |
| Ms. Pacman | 1249.79 | 1652.37 | **1861.80** | 1784.75 | 1310.62 | 1849.00 | 1920.25 | 2311.00 |
| Name This Game | 6960.46 | 7075.53 | **7560.33** | 5757.03 | 6094.08 | 7358.25 | 7565.67 | 7257.00 |
| Pong | 5.53 | **16.49** | -2.68 | 12.83 | 8.65 | 2.60 | 20.23 | 18.90 |
| Private Eye | 471.76 | **3609.96** | 1388.45 | 269.28 | 714.97 | 1277.53 | 7940.27 | 1788.00 |
| Q*Bert | 785.00 | 1074.77 | 2037.21 | 1215.42 | 3192.08 | **3955.10** | 2437.83 | 10596.00 |
| River Raid | 3460.62 | 4268.28 | 3676.12 | 4178.92 | **6005.62** | 4643.62 | 5671.51 | 8316.00 |
| Road Runner | 10086.74 | 15681.49 | 8978.17 | 17137.92 | 9390.83 | **19081.55** | 28286.88 | 18257.00 |
| Robotank | 11.65 | 15.34 | **16.11** | 6.46 | 9.90 | 12.17 | 20.73 | 51.60 |
| Seaquest | 1380.67 | 1926.10 | 762.10 | 1955.67 | 2275.83 | **2710.33** | 5313.43 | 5286.00 |
| Space Invaders | 797.29 | **1058.25** | 755.95 | 762.54 | 783.35 | 869.83 | 1148.21 | 1976.00 |
| Star Gunner | 2737.08 | **3892.51** | 708.66 | 2629.17 | 2856.67 | 1710.83 | 17462.88 | 57997.00 |
| Tennis | -3.41 | -0.96 | **0.00** | -10.32 | -2.50 | -6.37 | -0.93 | -2.50 |
| Time Pilot | 3505.42 | **4567.18** | 3076.98 | 4434.17 | 3651.25 | 4012.50 | 4567.18 | 5947.00 |
| Tutankham | 204.83 | 239.51 | 165.27 | 255.74 | 156.16 | **247.81** | 299.11 | 186.70 |
| Up and Down | 6841.83 | 6754.11 | **9468.04** | 7397.29 | 7574.53 | 6706.83 | 10984.70 | 8456.00 |
| Venture | 105.10 | **194.89** | 96.70 | 60.40 | 50.85 | 106.67 | 242.56 | 380.00 |
| Video Pinball | **84859.24** | 78405.27 | 17803.69 | 55646.66 | 18346.58 | 38528.58 | 84695.96 | 42684.00 |
| Wizard of Wor | 1249.89 | **2030.63** | 529.85 | 1175.24 | 1083.69 | 1177.08 | 4185.40 | 3393.00 |
| Zaxxon | 3221.67 | 3487.38 | 685.84 | **3928.33** | 596.67 | 2467.92 | 6548.52 | 4977.00 |

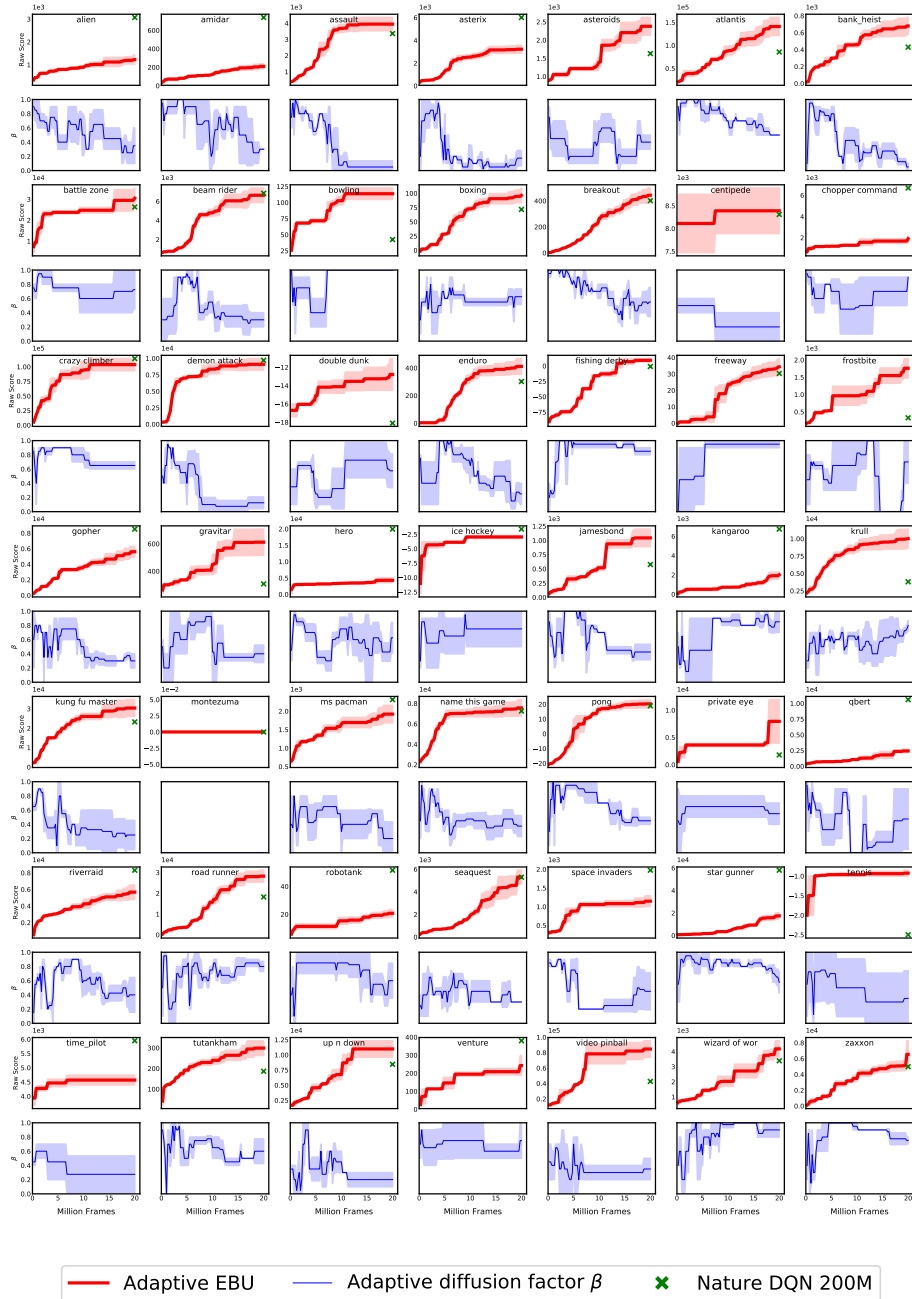# Appendix C    Learning curves and corresponding adaptive diffusion factor



Figure 7: Test scores and diffusion factor of Adaptive EBU. We report the mean and the standard deviation from 4 random seeds. We compare the performance of adaptive EBU with the result reported in Nature DQN, trained for 200M frames. The blue curve below each test score plot shows how adaptive EBU adapts its diffusion factor during the course of training.

# Appendix D    Network structure and hyperparameters

## 2D MNIST Maze Environment

Each state is given as a grey scale $28 \times 28$ image. We apply 2 convolutional neural networks (CNNs) and one fully connected layer to get the output $Q$-values for 4 actions: up, down, left and right. The first CNN uses 64 channels with $4 \times 4$ kernels and stride of 3. The next CNN uses 64 channels with $3 \times 3$ kernels and stride of 1. Then the layer is fully connected into a size of 512. Then we fully connect the layer into a size of the action space 4. After each layer, we apply a rectified linear unit.

We train the agent for a total of 200,000 steps. The agent performs $\epsilon$-greedy exploration. $\epsilon$ starts from 1 and is annealed to 0 at 200,000 steps in a quadratic manner: $\epsilon = \frac{1}{(200,000)^2}(\text{step} - 200,000)^2$. We use RMSProp optimizer with a learning rate of 0.001. The online-network is updated every 50 steps, the target network is updated every 2000 steps. The replay memory size is 30000 and we use minibatch size of 350. We use a discount factor $\gamma = 0.9$ and a diffusion factor $\beta = 1.0$. The agent plays the game until it reaches the goal or it stays in the maze for more than 1000 time steps.

## 49 Games of Atari 2600 Domain

**Common specifications for all baselines**
Almost all specifications such as hyperparameters and network structures are identical for all baselines. We use exactly the same network structure and hyperparameters of Nature DQN (Mnih et al., 2015). The raw observation is preprocessed into a gray scale image of $84 \times 84$. Then it passes through three convolutional layers: 32 channels with $8 \times 8$ kernels with a stride of 4; 64 channels with $4 \times 4$ kernels with a stride of 2; 64 channels with $3 \times 3$ kernels with a stride of 1. Then it is fully connected into a size of 512. Then it is again fully connected into the size of the action space.

We train baselines for 10M frames each, which is equivalent to 2.5M steps with frameskip of 4. The agent performs $\epsilon$-greedy exploration. $\epsilon$ starts from 1 and is linearly annealed to reach the final value 0.1 at 4M frames of training. We adopt 30 no-op evaluation methods. We use 8 random seeds for 10M frames and 4 random seeds for 20M frames. The network is trained by RMSProp optimizer with a learning rate of 0.00025. At each update (4 agent steps or 16 frames), we update transitions in minibatch with size 32. The replay buffer size is 1 million steps (4M frames). The target network is updated every 10,000 steps. The discount factor is $\gamma = 0.99$.

We divide the training process into 40 epochs (80 epochs for 20M frames) of 250,000 frames each. At the end of each epoch, the agent is tested for 30 episodes with $\epsilon = 0.05$. The agent plays the game until it runs out of lives or time (18,000 frames, 5 minutes in real time).

Below are detailed specifications for each algorithm.

**1. Episodic Backward Update**
We used $\beta = 0.5$ for the version EBU with constant diffusion factor. For adaptive EBU, we used 11 parallel learners ($K = 11$) with diffusion factors 0.0, 0.1, ..., 1.0. We synchronize the learners at every 250,000 frames ($B = 62,500$ steps).

**2. Prioritized Experience Replay**
We use the rank-based DQN version of Prioritized ER and use the hyperparameters chosen by the authors (Schaul et al., 2016): $\alpha = 0.5 \rightarrow 0$ and $\beta = 0$.

**3. Retrace($\lambda$)**
Just as EBU, we sample a random episode and then generate the Retrace target for the transitions in the sampled episode. We follow the same evaluation process as that of Munos et al., 2016. First, we calculate the trace coefficients from $s = 1$ to $s = T$ (terminal).

$$c_s = \lambda \min \left(1, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)}\right) \tag{3}$$

Where $\mu$ is the behavior policy of the sampled transition and the evaluation policy $\pi$ is the current policy. Then we generate a loss vector for transitions in the sample episode from $t = T$ to $t = 1$.

$$\Delta Q(x_{t-1}, a_{t-1}) = c_t \lambda \Delta Q(x_t, a_t) + [r(x_{t-1}, a_{t-1}) + \gamma \mathbb{E}_\pi Q(x_t, :) - Q(x_{t-1}, a_{t-1})]. \tag{4}$$

**4. Optimality Tightening**
We use the source code (`https://github.com/ShibiHe/Q-Optimality-Tightening`), modify the maximum test steps and test score calculation to match the evaluation policy of Nature DQN.

# Appendix E      Supplementary figure: backward update algorithm

Line #7 of Algorithm 2: Sample a random episode $E$.



Line # 8~9: Generate a temporary target Q table $\tilde{Q}$ with the next state vector $S'$. Initialize a target vector $y$.
Let there be $n$ possible actions in the environment. $\mathcal{A} = \{a^{(1)}, a^{(2)}, \dots, a^{(n)}\}$.
Note that $\hat{Q}$ is the target Q-value and $\hat{Q}(S_{T+1}, :) = 0$.



Line # 10~12, first iteration (k = T-1): Update $\tilde{Q}$ and $y$. Let the T-th action in the replay memory be $A_T = a^{(2)}$.

     ① line # 15: update $\tilde{Q}[A_{k+1}, k] = \tilde{Q}[A_T, T-1] = \tilde{Q}[a^{(2)}, T-1] \leftarrow \beta y_T + (1-\beta)\hat{Q}(S_T, a^{(2)})$

     ② line # 16: update $y_k = y_{T-1} \leftarrow R_{T-1} + \gamma \max \tilde{Q}[:, T-1]$



Line # 10~12, second iteration (k = T-2): Update $\tilde{Q}$ and $y$. Let the (T-1)-th action in the replay memory be $A_{T-1} = a^{(1)}$.

     ① line # 15: update $\tilde{Q}[A_{k+1}, k] = \tilde{Q}[A_{T-1}, T-2] = \tilde{Q}[a^{(1)}, T-2] \leftarrow \beta y_{T-1} + (1-\beta)\hat{Q}(S_{T-1}, a^{(1)})$

     ② line # 16: update $y_k = y_{T-2} \leftarrow R_{T-2} + \gamma \max \tilde{Q}[:, T-2]$



Repeat this update until k =1.

Figure 8: Target generation process from the sampled episode E

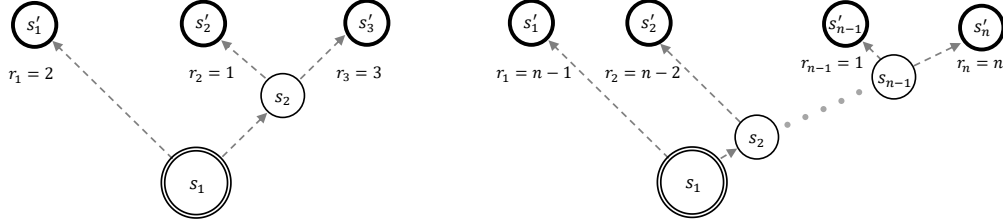## Appendix F        Comparison to other multi-step methods.



Figure 9: A motivating example where $Q(\lambda)$ underperforms Episodic Backward Update. **Left:** A simple navigation domain with 3 possible episodes. $s_1$ is the initial state. States with ' signs are the terminal states. **Right:** An extended example with $n$ possible episodes.

Imagine a toy navigation environment as in Figure 9, left. Assume that an agent has experienced all possible trajectories: $(s_1 \rightarrow s'_1)$; $(s_1 \rightarrow s_2 \rightarrow s'_2)$ and $(s_1 \rightarrow s_2 \rightarrow s'_3)$. Let the discount factor $\gamma$ be 1. Then optimal policy is $(s_1 \rightarrow s_2 \rightarrow s'_3)$. With a slight abuse of notation let $Q(s_i, s_j)$ denote the value of the action that leads to the state $s_j$ from the state $s_i$. We will show that $Q(\lambda)$ and $Q^*(\lambda)$ methods underperform Episodic Backward Update in such examples with many suboptimal branching paths.

$Q(\lambda)$ method cuts trace of the path when the path does not follow greedy actions given the current $Q$-value. For example, assume a $Q(\lambda)$ agent has updated the value $Q(s_1, s'_1)$ at first. When the agent tries to update the values of the episode $(s_1 \rightarrow s_2 \rightarrow s'_3)$, the greedy policy of the state $s_1$ heads to $s'_1$. Therefore the trace of the optimal path is cut and the reward signal $r_3$ is not passed to $Q(s_1, s_2)$. This problem becomes more severe if the number of suboptimal branches increases as illustrated in Figure 9, right. Other variants of $Q(\lambda)$ algorithm that cut traces, such as Retrace($\lambda$), have the same problem. EBU does not suffer from this issue, because EBU does not cut the trace, but performs max operations at every branch to propagate the maximum value.

$Q^*(\lambda)$ is free from the issues mentioned above since it does not cut traces. However, to guarantee convergence to the optimal value function, it requires the parameter $\lambda$ to be less than $\frac{1-\gamma}{2\gamma}$. In convention, the discount factor $\gamma \approx 1$. For a small value of $\lambda$ that satisfies the constraint, the update of distant returns becomes nearly negligible. However, EBU does not have any constraint of the diffusion factor $\beta$ to guarantee convergence.

# Appendix G    Theoretical guarantees

Now, we will prove that the episodic backward update algorithm converges to the true action-value function $Q^*$ in the case of finite and deterministic environment.

**Definition 1.** *(Deterministic MDP)*

$M = (\mathcal{S}, \mathcal{A}, P, R)$ *is a **deterministic MDP** if $\exists g : \mathcal{S} \times \mathcal{A} \to \mathcal{S}$ s.t.*

$$P(s'|s, a) = \begin{cases} 1 & \text{if } s' = g(s, a) \\ 0 & \text{else} \end{cases} \quad \forall (s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S},$$

In the episodic backward update algorithm, a single (state, action) pair can be updated through multiple episodes, where the evaluated targets of each episode can be different from each other. Therefore, unlike the bellman operator, episodic backward operator depends on the exploration policy for the MDP. Therefore, instead of expressing different policies in each state, we define a schedule to represent the frequency of every distinct episode (which terminates or continues indefinitely) starting from the target (state, action) pair.

**Definition 2.** *(Schedule)*

*Assume a MDP $M = (\mathcal{S}, \mathcal{A}, P, R)$, where $R$ is a bounded function. Then, for each state $(s, a) \in \mathcal{S} \times \mathcal{A}$ and $j \in [1, \infty]$, we define $j$-**length path set** $p_{s,a}(j)$ and **path set** $p(s, a)$ for $(s, a)$ as*

$$p_{s,a}(j) = \left\{ (s_i, a_i)_{i=0}^{j} | (s_0, a_0) = (s, a), P(s_{i+1}|s_i, a_i) > 0 \quad \forall i \in [0, j-1], s_j \quad is \quad terminal \right\}.$$

*and $p_{s,a} = \cup_{j=1}^{\infty} p_{s,a}(j)$.*

*Also, we define a **schedule set** $\lambda_{s,a}$ for (state action) pair $(s, a)$ as*

$$\lambda_{s,a} = \left\{ (\lambda_i)_{i=1}^{|p_{s,a}|} | \sum_{i=1}^{|p_{s,a}|} \lambda_i = 1, \lambda_i > 0 \quad \forall i \in [1, |p_{s,a}|] \right\}.$$

*Finally, to express the varying schedule in time at the RL scenario, we define a **time schedule set** $\lambda$ for MDP $M$ as*

$$\lambda = \left\{ \{\lambda_{s,a}(t)\}_{(s,a) \in \mathcal{S} \times \mathcal{A}, t=1}^{\infty} | \lambda_{s,a}(t) \in \lambda_{s,a}, \forall (s, a) \in \mathcal{S} \times \mathcal{A}, t \in [1, \infty] \right\}.$$

Since no element of the path can be the prefix of the others, the path set corresponds to the enumeration of all possible episodes starting from each (state, action) pair. Therefore, if we utilize multiple episodes from any given policy, we can see the empirical frequency for each path in the path set belongs to the schedule set. Finally, since the exploration policy can vary across time, we can group independent schedules into the time schedule set.

For a given time schedule and MDP, now we define the episodic backward operator.

**Definition 3.** *(Episodic backward operator)*

*For an MDP $M = (\mathcal{S}, \mathcal{A}, P, R)$, and a time schedule $\{\lambda_{s,a}(t)\}_{t=1,(s,a) \in \mathcal{S} \times \mathcal{A}}^{\infty} \in \lambda$.*

*Then, the **episodic backward operator** $H_t^{\beta}$ is defined as*

$$(H_t^{\beta} Q)(s, a) \tag{5}$$

$$= \mathbb{E}_{s' \in \mathcal{S}, P(s'|s,a)} \left[ r(s, a, s') + \gamma \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)}(t))_i \mathbb{1}(s_{i1} = s') \left[ \max_{1 \leq j \leq |(p_{s,a})_i|} T_{(p_{s,a})_i}^{\beta, Q}(j) \right] \right].$$

$$T_{(p_{s,a})_i}^{\beta, Q}(j) \tag{6}$$

$$= \sum_{k=1}^{j-1} \beta^{k-1} \gamma^{k-1} \left\{ \beta r(s_{ik}, a_{ik}, s_{i(k+1)}) + (1 - \beta) Q(s_{ik}, a_{ik}) \right\} + \beta^{j-1} \gamma^{j-1} \max_{a \neq a_j} Q(s_{ij}, a_{ij}).$$

*Where $(p_{s,a})_i$ is the $i$-th path of the path set, and $(s_{ij}, a_{ij})$ corresponds to the $j$-th (state, action) pair of the $i$-th path.*

Episodic backward operator consists of two parts. First, given the path that initiates from the target (state, action) pair, the function $T^{\beta,Q}_{(p_{s,a})_i}$ computes the maximum return of the path via backward update. Then, the return is averaged by every path in the path set. Now, if the MDP $M$ is deterministic, we can prove that the episodic backward operator is a contraction in the sup-norm, and the fixed point of the episodic backward operator is the optimal action-value function of the MDP regardless of the time schedule.

**Theorem 3.** *(Contraction of the episodic backward operator and the fixed point)*

*Suppose $M = (\mathcal{S}, \mathcal{A}, P, R)$ is a deterministic MDP. Then, for any time schedule $\{\lambda_{s,a}(t)\}^{\infty}_{t=1,(s,a)\in\mathcal{S}\times\mathcal{A}} \in \lambda$, $H^{\beta}_t$ is a contraction in the sup-norm for any t, i.e*

$$\|(H^{\beta}_t Q_1) - (H^{\beta}_t Q_2)\|_{\infty} \le \gamma\|Q_1 - Q_2\|_{\infty}. \tag{7}$$

*Furthermore, for any time schedule $\{\lambda_{s,a}(t)\}^{\infty}_{t=1,(s,a)\in\mathcal{S}\times\mathcal{A}} \in \lambda$, the fixed point of $H^{\beta}_t$ is the optimal Q function $Q^*$.*

*Proof.* First, we prove $T^{\beta,Q}_{(p_{s,a})_i}(j)$ is a contraction in the sup-norm for all $j$.

Since $M$ is a deterministic MDP, we can reduce the return as

$$T^{\beta,Q}_{(p_{s,a})_i}(j) = \left( \sum_{k=1}^{j-1} \beta^{k-1}\gamma^{k-1} \{\beta r(s_{ik}, a_{ik}) + (1-\beta)Q(s_{ik}, a_{ik})\} + \beta^{j-1}\gamma^{j-1} \max_{a\ne a_j} Q(s_{ij}, a_{ij}) \right). \tag{8}$$

$$\begin{aligned}
\|T^{\beta,Q_1}_{(p_{s,a})_i}(j) - T^{\beta,Q_2}_{(p_{s,a})_i}(j)\|_{\infty} &\le \left\{ (1-\beta) \sum_{k=1}^{j-1} \beta^{k-1}\gamma^{k-1} + \beta^{j-1}\gamma^{j-1} \right\} \|Q_1 - Q_2\|_{\infty} \\
&= \left\{ \frac{(1-\beta)(1-(\beta\gamma)^{j-1})}{1-\beta\gamma} + \beta^{j-1}\gamma^{j-1} \right\} \|Q_1 - Q_2\|_{\infty} \\
&= \frac{1 - \beta + \beta^j\gamma^{j-1} - \beta^j\gamma^j}{1-\beta\gamma} \|Q_1 - Q_2\|_{\infty} \\
&= \left\{ 1 + (1-\gamma) \frac{\beta^j\gamma^{j-1} - \beta}{1-\beta\gamma} \right\} \|Q_1 - Q_2\|_{\infty} \\
&\le \|Q_1 - Q_2\|_{\infty} \quad (\because \beta \in [0,1], \gamma \in [0,1)).
\end{aligned} \tag{9}$$

Also, at the deterministic MDP, the episodic backward operator can be reduced to

$$(H^{\beta}_t Q)(s, a) = r(s, a) + \gamma \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)})_i(t) \left[ \max_{1 \le j \le |(p_{s,a})_i|} T^{\beta,Q}_{(p_{s,a})_i}(j) \right]. \tag{10}$$

18

Therefore, we can finally conclude that

$$\|(H_t^\beta Q_1) - (H_t^\beta Q_2)\|_\infty$$

$$= \max_{s,a} \left| H_t^\beta Q_1(s,a) - H_t^\beta Q_2(s,a) \right|$$

$$\leq \gamma \max_{s,a} \left[ \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)}(t))_i \left| \left\{ \max_{1 \leq j \leq |(p_{s,a})_i|} T_{(p_{s,a})_i}^{\beta,Q_1}(j) \right\} - \left\{ \max_{1 \leq j \leq |(p_{s,a})_i|} T_{(p_{s,a})_i}^{\beta,Q_2}(j) \right\} \right| \right]$$

$$\leq \gamma \max_{s,a} \left[ \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)}(t))_i \max_{1 \leq j \leq |(p_{s,a})_i|} \left\{ \left| T_{(p_{s,a})_i}^{\beta,Q_1}(j) - T_{(p_{s,a})_i}^{\beta,Q_2}(j) \right| \right\} \right]$$

$$\leq \gamma \max_{s,a} \left[ \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)}(t))_i \|Q_1 - Q_2\|_\infty \right]$$

$$= \gamma \max_{s,a} \left[ \|Q_1 - Q_2\|_\infty \right]$$

$$= \gamma \|Q_1 - Q_2\|_\infty. \tag{11}$$

Therefore, we have proved that the episodic backward operator is a contraction independent of the schedule. Finally, we prove that the distinct episodic backward operators in terms of schedule have the same fixed point, $Q^*$. A sufficient condition to prove this is given by

$$\left[ \max_{1 \leq j \leq |(p_{s,a})_i|} T_{(p_{s,a})_i}^{\beta,Q^*}(j) \right] = \frac{Q^*(s,a) - r(s,a)}{\gamma} \; \forall 1 \leq i \leq |p_{s,a}|.$$

We will prove this by contradiction. Assume $\exists i$ s.t. $\left[ \max_{1 \leq j \leq |(p_{s,a})_i|} T_{(p_{s,a})_i}^{\beta,Q^*}(j) \right] \neq \frac{Q^*(s,a) - r(s,a)}{\gamma}$.

First, by the definition of $Q^*$ fuction, we can bound $Q^*(s_{ik}, a_{ik})$ and $Q^*(s_{ik}, :)$ for every $k \geq 1$ as follows.

$$Q^*(s_{ik}, a) \leq \gamma^{-k} Q^*(s,a) - \sum_{m=0}^{k-1} \gamma^{m-k} r(s_{im}, a_{im}). \tag{12}$$

Note that the equality holds if and only if the path $(s_i, a_i)_{i=0}^{k-1}$ is the optimal path among the ones that start from $(s_0, a_0)$. Therefore, $\forall 1 \leq j \leq \left| (p_{s,a})_i \right|$, we can bound $T_{(p_{s,a})_i}^{\beta,Q^*}(j)$.

$$T^{\beta,Q}_{(p_{s,a})_i}(j)$$

$$= \sum_{k=1}^{j-1} \beta^{k-1}\gamma^{k-1} \{\beta r(s_{ik}, a_{ik}) + (1-\beta)Q(s_{ik}, a_{ik})\} + \beta^{j-1}\gamma^{j-1} \max_{a \neq a_j} Q(s_{ij}, a_{ij})$$

$$\leq \left\{ (\sum_{k=1}^{j-1}(1-\beta)\beta^{k-1}) + \beta^{j-1} \right\} \gamma^{-1}Q^*(s,a)$$

$$+ \sum_{k=1}^{j-1} \left\{ \beta^{k-1}\gamma^{k-1} \left( \beta r(s_{ik}, a_{ik}) - \sum_{m=0}^{k-1}(1-\beta)\gamma^{m-k}r(s_{im}, a_{im}) \right) \right\}$$

$$- \sum_{m=0}^{j-1} \beta^{j-1}\gamma^{j-1}\gamma^{m-j}r(s_{im}, a_{im})$$

$$= \gamma^{-1}Q^*(s,a) + \sum_{k=1}^{j-1} \beta^k \gamma^{k-1} r(s_{ik}, a_{ik})$$

$$- \sum_{m=0}^{j-2} \left\{ \sum_{k=m+1}^{j-1}(1-\beta)\beta^{k-1}\gamma^{m-1}r(s_{im}, a_{im}) \right\} - \sum_{m=0}^{j-1} \beta^{j-1}\gamma^{m-1}r(s_{im}, a_{im})$$

$$= \gamma^{-1}Q^*(s,a) + \sum_{m=1}^{j-1} \beta^m \gamma^{m-1} r(s_{im}, a_{im})$$

$$- \sum_{m=0}^{j-2}(\beta^m - \beta^{j-1})\gamma^{m-1}r(s_{im}, a_{im}) - \sum_{m=0}^{j-1} \beta^{j-1}\gamma^{m-1}r(s_{im}, a_{im})$$

$$= \gamma^{-1}Q^*(s,a) - \gamma^{-1}r(s_{i0}, a_{i0}) = \frac{Q^*(s,a) - r(s,a)}{\gamma}.$$

(13)

Since this occurs for any arbitrary path, the only remaining case is when

$$\exists i \text{ s.t. } \left[ \max_{1 \leq j \leq |(p_{s,a})_i|} T^{\beta,Q^*}_{(p_{s,a})_i}(j) \right] < \frac{Q^*(s,a) - r(s,a)}{\gamma}.$$

Now, let's turn our attention to the path $s_0, s_1, s_2, ...., s_{|(p_{s,a})_i|}$. Let's first prove the contradiction when the length of the contradictory path is finite. If $Q^*(s_{i1}, a_{i1}) < \gamma^{-1}(Q^*(s,a) - r(s,a))$, then by the Bellman equation, there exists an action $a \neq a_{i1}$ s.t. $Q^*(s_{i1}, a) = \gamma^{-1}(Q^*(s,a) - r(s,a))$. Then, we can find that $T^{\beta,Q^*}_{(p_{s,a})_1}(1) = \gamma^{-1}(Q^*(s,a) - r(s,a))$. It contradicts the assumption, therefore $a_{i1}$ should be the optimal action in $s_{i1}$.

Repeating the procedure, we conclude that $a_{i1}, a_{i2}, ..., a_{|(p_{s,a})_i|-1}$ are optimal with respect to their corresponding states.

Finally, $T^{\beta,Q^*}_{(p_{s,a})_1}(|(p_{s,a})_i|) = \gamma^{-1}(Q^*(s,a) - r(s,a))$ since all the actions satisfy the optimality condition of the inequality in equation 7. Therefore, it contradicts the assumption.

In the case of an infinite path, we will prove that for any $\epsilon > 0$, there is no path that satisfies $\frac{Q^*(s,a) - r(s,a)}{\gamma} - \left[ \max_{1 \leq j \leq |(p_{s,a})_i|} T^{\beta,Q^*}_{(p_{s,a})_i}(j) \right] = \epsilon.$

Since the reward function is bounded, we can define $r_{\max}$ as the supremum norm of the reward function. Define $q_{\max} = \max_{s,a} |Q(s,a)|$ and $R_{\max} = \max\{r_{\max}, q_{\max}\}$. We can assume $R_{\max} > 0$. Then, let's set $n_\epsilon = \lceil \log_\gamma \frac{\epsilon(1-\gamma)}{R_{\max}} \rceil + 1$. Since $\gamma \in [0,1)$, $R_{\max} \frac{\gamma^{n_\epsilon}}{1-\gamma} < \epsilon$. Therefore, by applying the procedure on the finite path case for $1 \le j \le n_\epsilon$, we can conclude that the assumption leads to a contradiction. Since the previous $n_\epsilon$ trajectories are optimal, the rest trajectories can only generate a return less than $\epsilon$.

Finally, we proved that $\left[\max_{1 \le j \le |(p_{s,a})_i|} T^{\beta,Q^*}_{(p_{s,a})_i}(j)\right] = \frac{Q^*(s,a) - r(s,a)}{\gamma}$ $\forall 1 \le i \le |p_{s,a}|$ and therefore, every episodic backward operator has $Q^*$ as the fixed point. □

Finally, we will show that the online episodic backward update algorithm converges to the optimal $Q$ function $Q^*$.

**Restatement of Theorem 1.** *Given a finite, deterministic, and tabular MDP $M = (\mathcal{S}, \mathcal{A}, P, R)$, the episodic backward update algorithm, given by the update rule*

$Q_{t+1}(s_t, a_t)$

$= (1-\alpha_t)Q_t(s_t, a_t) + \alpha_t \left[ r(s_t, a_t) + \gamma \sum_{i=1}^{|p_{s_t, a_t}|} (\lambda_{(s_t, a_t)})_i(t) \left[\max_{1 \le j \le |(p_{s_t, a_t})_i|} T^{\beta,Q}_{(p_{s_t, a_t})_i}(j)\right] \right]$

*converges to the optimal Q-function w.p. 1 as long as*

- *The step size satisfies the Robbins-Monro condition;*

- *The sample trajectories are finite in lengths $l$: $\mathbb{E}[l] < \infty$;*

- *Every (state, action) pair is visited infinitely often.*

For the proof of Theorem 1, we follow the proof of Melo, 2001.

**Lemma 1.** *The random process $\Delta_t$ taking values in $\mathbb{R}^n$ and defined as*

$\Delta_{t+1}(x) = (1 - \alpha_t(x))\Delta_t(x) + \alpha_t(x)F_t(x)$

*converges to zero w.p. 1 under the following assumptions:*

- $0 \le \alpha_t \le 1, \sum_t \alpha_t(x) = \infty$ *and* $\sum_t \alpha_t^2(x) < \infty$;

- $\|\mathbb{E}[F_t(x)|\mathcal{F}_t]\|_W \le \gamma \|\Delta_t\|_W$, *with $\gamma < 1$;*

- $\mathbf{var}[F_t(x)|\mathcal{F}_t] \le C \left(1 + \|\Delta_t\|_W^2\right)$, *for $C > 0$.*

By Lemma 1, we can prove that the online episodic backward update algorithm converges to the optimal $Q^*$.

*Proof.* First, by assumption, the first condition of Lemma 1 is satisfied. Also, we can see that by substituting $\Delta_t(s,a) = Q_t(s,a) - Q^*(s,a)$, and $F_t(s,a) = r(s,a) + \gamma \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)})_i(t) \left[\max_{1 \le j \le |(p_{s,a})_i|} T^{\beta,Q}_{(p_{s,a})_i}(j)\right] - Q^*(s,a)$. $\|\mathbb{E}[F_t(s,a)|\mathcal{F}_t]\|_\infty = \|(H_t^\beta Q_t)(s,a) - (H_t^\beta Q^*)(s,a)\|_\infty \le \gamma \|\Delta_t\|_\infty$, where the inequality holds due to the contraction of the episodic backward operator.

Then, $\mathbf{var}[F_t(x)|\mathcal{F}_t] = \mathbf{var}\left[r(s,a) + \gamma \sum_{i=1}^{|p_{s,a}|} (\lambda_{(s,a)})_i(t) \left[\max_{1 \le j \le |(p_{s,a})_i|} T^{\beta,Q}_{(p_{s,a})_i}(j)\right] \Big| \mathcal{F}_t\right]$.

Since the reward function is bounded, the third condition also holds as well. Finally, by Lemma 1, $Q_t$ converges to $Q^*$.

□

Although the episodic backward operator can accommodate infinite paths, the operator can be practical when the maximum length of the episode is finite. This assumption holds for many RL domains, such as the ALE.