Reinforcement Learning with Adaptive Temporal Discounting

Anonymous authors Paper under double-blind review

Keywords: Adaptive Temporal Discounting, Weber-Fechner law, Log-compressed Timeline.

Summary

Conventional reinforcement learning (RL) methods often fix a single discount factor for future rewards, limiting their ability to handle diverse temporal requirements. We propose a framework that trains an agent across a spectrum of discount factors—interpreting each value function as a sample of a Laplace transform—then applies an inverse transform to recover a log-compressed representation of expected future reward. This representation enables post hoc adjustments to the discount function (e.g., exponential, hyperbolic, or finite horizon) without retraining. Furthermore, by precomputing a library of policies, the agent can dynamically select which policy maximizes a newly specified discount objective at runtime, effectively constructing a hybrid policy in environments with shifting deadlines or reward structures. The log-compressed timeline aligns with human temporal perception, as described by the Weber-Fechner law, maintaining uniform relative precision across timescales thus enhancing efficiency in scale-free environments. We demonstrate this framework in a grid-world navigation task, where the agent adapts to varying time horizons.

Contribution(s)

 We describe a method for computing log-compressed representation of expected future reward and demonstrate that when a library of policies is available, this representation enables immediate re-evaluation of these policies under any desired discount function.

Context: Prior work established methods for computing log-compressed representation of expected future reward (Momennejad & Howard, 2018; Tiganj et al., 2019; Tano et al., 2020; Masset et al., 2023). We extend prior work to dynamical policy evaluation setting with arbitrary temporal discounting.

2. We demonstrate that a logarithmic representation of future time aligns with human temporal perception (Weber-Fechner law) and provides efficient decision-making in scale-free environments by maintaining uniform relative precision across time scales.

Context: The Weber-Fechner law Fechner (1860/1912) is a widely referenced principle in psychophysics, stating that perceived magnitude is proportional to the logarithm of stimulus intensity, implying a logarithmic scale.

Reinforcement Learning with Adaptive Temporal Discounting

Anonymous authors

Paper under double-blind review

Abstract

| 1 | Conventional reinforcement learning (RL) methods often fix a single discount fac- |
|----|---|
| 2 | tor for future rewards, limiting their ability to handle diverse temporal requirements. |
| 3 | We propose a framework that trains an agent across a spectrum of discount fac- |
| 4 | tors-interpreting each value function as a sample of a Laplace transform-then applies |
| 5 | an inverse transform to recover a log-compressed representation of expected future re- |
| 6 | ward. This representation enables post hoc adjustments to the discount function (e.g., |
| 7 | exponential, hyperbolic, or finite horizon) without retraining. Furthermore, by precom- |
| 8 | puting a library of policies, the agent can dynamically select which policy maximizes |
| 9 | a newly specified discount objective at runtime, effectively constructing a hybrid pol- |
| 10 | icy in environments with shifting deadlines or reward structures. The log-compressed |
| 11 | timeline aligns with human temporal perception, as described by the Weber-Fechner |
| 12 | law, maintaining uniform relative precision across timescales thus enhancing efficiency |
| 13 | in scale-free environments. We demonstrate this framework in a grid-world navigation |
| 14 | task, where the agent adapts to varying time horizons. |

15 1 Introduction

16 In traditional RL setting, agents learn by maximizing a cumulative future reward, which is typi-17 cally exponentially discounted over time using a fixed temporal discount factor (γ) (Sutton & Barto, 18 2018). Despite its widespread adoption, such discounting approach presents limitations, particularly 19 in dynamic environments where the relevance of future rewards may vary over time. For example, 20 if an agent encounters an emergency situation where immediate action is critical a high discount on 21 future rewards is beneficial. Conversely, in investment scenarios, considering longer-term outcomes 22 might be preferred in some cases, necessitating a lower discount rate while in other cases the focus 23 might be on short-term investments where a higher discount rate would be preferred. Furthermore, 24 while the exponential discounting has been widely used since exponentially discounted future re-25 ward can be efficiently computed using temporal difference (TD) learning thanks to the efficiency 26 of Bellman equation, other types of temporal discounting are often preferred. Specifically, hyper-27 bolic discounting has gained attention for its ability to more closely mimic human decision-making 28 processes. Unlike exponential discounting, hyperbolic discounting reduces the value of rewards less 29 steeply over time, which has been shown to reflect more accurately how people value immediate 30 versus delayed rewards (Ainslie, 1975). This form of discounting suggests that people often display 31 a preference for smaller-sooner rewards over larger-later rewards when the delays are short, but this 32 preference can switch as the delays increase, a phenomenon often referred to as "temporal inconsis-33 tency" or "preference reversal" (Green & Myerson, 1994; Laibson, 1997). More generally, humans 34 can adjust the discounting function to adapt to the temporal statistics of the environment (Redish, 35 2004; Frederick et al., 2002). These considerations highlight the need for more flexible approaches 36 to temporal discounting in RL, where discounting can be adjusted dynamically after the agent has 37 been trained.

38 1.1 Related work

39 Limitations of exponential discounting were addressed in a number of previous studies. Inspired

40 by human decision-making, previous work mainly focused on hyperbolic discounting by integrating

41 over a spectrum of γ values (Kurth-Nelson & Redish, 2009; Fedus et al., 2019; Masset et al., 2023;

42 Tiganj et al., 2019; Tano et al., 2020). Tang et al. (2021) used Taylor expansion of discount rates to

43 interpolate value functions of two distinct discount factors.

44 Going beyond integration across γ values, previous work proposed the use of inverse Laplace trans-45 form to convert a value function across the set of γ values into a function of future time (Momenne-46 jad & Howard, 2018; Tiganj et al., 2019; Tano et al., 2020; Masset et al., 2023). This was primarily 47 applied to modeling of human decision making as existence of a mental timeline of the future has been supported by behavioral (Tiganj et al., 2022) and neural evidence (Cao et al., 2024). Proposed 48 49 approach builds on this work and integrates the Laplace framework into a setting where agents need 50 to select actions at every step. We propose that agents construct multiple timelines of the future and 51 use policy selection at every step. 52

This approach is in general compatible with model-free and model-based setting. For example, Tano et al. (2020) used spectrum of γ values to compute a reward as a function of future time in a modelfree setting with TD learning. Model-based RL learning of arbitrary discount functions, including hyperbolic was also done using function approximation of the optimal policy (Schultheis et al., 2022). Model based RL approaches used successor representation and statistical learning computed with a spectrum of γ values to estimate timeline of the future (Tano et al., 2020; Momennejad & Howard, 2018; Tiganj et al., 2019).

59 **1.2** Summary of the proposed approach

Here we describe an RL method with adaptive temporal discounting, where the discounting function can be chosen after training and dynamically adjusted. This is achieved by computing expected reward as a function of future time τ^* for a policy π . Such function then enables temporal discounting with arbitrary functional forms, such as exponential and power-law as well as temporal windows (e.g., expected reward from time t_i to time t_j at a given state) by varying the weight given to different parts of the future timeline.

To compute expected reward as a function of future time for a given state under a policy π , we compute a set of values V for a spectrum exponential discount rates γ s. Using a linear transformation this set of values can then be converted from a function of γ into a function of future time, τ^* . From a set of candidate policies π , this approach can then select a policy that maximizes a reward under the desired temporal discounting.

In the limit where the value V is computed for all possible γ values (i.e., γ is a continuous variable) we show that this approach can select an optimal policy from a set of precomputed policies under a user-specified temporal discounting function. For practical applications, we consider a discrete set of geometrically spaced γ values. This gives rise to a log-compressed representation of the future time where the temporal resolution gradually decreases as a function of distance from the current state.

77 Log-compressed time reflects the Weber-Fechner law which governs scale-invariance in human per-78 ception across a number of domains, including time: perceived magnitude is proportional to the 79 logarithm of the true magnitude of the stimulus (Fechner, 1860/1912). In time perception the scale-80 invariance is manifested in behavioral tasks such as interval reproduction, where variability of repro-81 duced interval by human participants is proportional to the duration of the interval. In other words, 82 timing precision is proportional to the length of the interval itself, a phenomenon known as the 83 scalar property (Gibbon, 1977; Wilkes, 2015). Such a logarithmically compressed representation of 84 time is argued to be advantageous when the environment lacks a single dominant timescale (i.e., it is 85 scale-free), so that events or correlations extend across multiple orders of magnitude in a self-similar 86 manner (Howard & Shankar, 2018; Shankar & Howard, 2013). We show that the log-compressed

87 representation of the future time described here, captures the scalar property and enables efficient

- decision making under the assumption that variability in reward time (or spatial position) scales with
- 89 its temporal (or spatial) distance.

90 2 Methods

- We consider a standard RL framework modeled as a Markov Decision Process (MDP), defined by the tuple $\langle S, A, P, R \rangle$, where:
- 93 S represents the state space, encompassing all possible states the agent may encounter,
- 94 \mathcal{A} denotes the action space, consisting of all actions available to the agent,
- 95 P : S × A × S → [0, 1] is the state transition probability function, specifying the probability
 96 P(s'|s, a) of transitioning to state s' ∈ S from state s ∈ S upon taking action a ∈ A,
- 97 R: S × A × S → R is the reward function, defining the immediate reward R(s, a, s') received
 98 when transitioning from state s to state s' via action a.
- 99 The agent's goal is to learn a policy $\pi(a|s)$, which maps states to a probability distribution over
- 100 actions, maximizing the expected cumulative reward over time under a specified temporal discount-101 ing scheme. In the traditional RL formulation, this is captured by the state-value function with an
- 102 exponential discount factor $\gamma \in [0, 1]$:

$$V_{\gamma}^{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{\tau^*=0}^{\infty} \gamma^{\tau^*} r_{t+\tau^*} \mid s_t = s \right], \quad \forall s \in \mathcal{S},$$
(1)

103 where $r_{t+\tau^*} = \mathcal{R}(s_{t+\tau^*}, a_{t+\tau^*}, s_{t+\tau^*+1})$ is the reward at time step $t + \tau^*, \tau^*$ denotes future time 104 steps relative to the current time t, and the expectation \mathbb{E}_{π} is taken over trajectories induced by the 105 policy π and the transition dynamics \mathcal{P} . The discount factor γ exponentially weights future rewards, 106 prioritizing immediate rewards when $\gamma < 1$. However, this formulation fixes the temporal preference 107 at training time, limiting adaptability to different discounting schemes or planning horizons post-108 training.

109 2.1 Computing Expected Rewards as a Function of Future Time

110 To enable this adaptability, we compute the state-value function $V_{\gamma}^{\pi}(s)$ across a discrete set of dis-111 count factors $\Gamma = \{\gamma_1, \gamma_2, \dots, \gamma_m\}$, where each $\gamma_i \in [0, 1]$. For a specific discount factor γ_i , the 112 value function is:

$$V_{\gamma_i}^{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{\tau^*=0}^{\infty} \gamma_i^{\tau^*} r_{t+\tau^*} \mid s_t = s \right].$$
(2)

113 Define $g(\tau^*) = \mathbb{E}_{\pi} [r_{t+\tau^*} | s_t = s]$, which represents the expected reward at future time step τ^* 114 under policy π , starting from state s at time t. This allows us to rewrite the value function as:

$$V_{\gamma_i}^{\pi}(s) = \sum_{\tau^*=0}^{\infty} \gamma_i^{\tau^*} g(\tau^*).$$
(3)

115 This expression resembles the discrete analog of the Laplace transform (or unilateral Z-transform)

116 of the sequence $g(\tau^*)$, with γ_i acting as the transform variable. Equivalently, letting $\sigma_i = -\ln(\gamma_i)$ 117 (so $\gamma_i = e^{-\sigma_i}$), we can express it as:

$$V_{\gamma_i}^{\pi}(s) = \sum_{\tau^*=0}^{\infty} e^{-\sigma_i \tau^*} g(\tau^*),$$
(4)

118 where $V_{\gamma_i}^{\pi}(s)$ can be interpreted as the Laplace transform of $g(\tau^*)$ evaluated at σ_i , i.e., $V_{\gamma_i}^{\pi}(s) = 119 \quad \mathcal{L}\{g(\tau^*)\}(\sigma_i)$. By computing $V_{\gamma_i}^{\pi}(s)$ for multiple $\gamma_i \in \Gamma$, we obtain a set of samples of this trans-120 form at points σ_i :

$$\mathbf{V}^{\pi}(s) = [V_{\gamma_1}^{\pi}(s), V_{\gamma_2}^{\pi}(s), \dots, V_{\gamma_m}^{\pi}(s)].$$

(5)

121 To recover the expected reward function $g(\tau^*)$ from these samples, we apply an inverse transform. 122 With a finite set of γ_i , we approximate $g(\tau^*)$ for discrete $\tau^* = 0, 1, 2, ...$ using numerical methods, 123 such as the Post inversion formula Post (1930) (see Sec. S1 for discussion on numerical stability 124 and demonstrations of other approaches for computing the inverse that can provide better numerical 125 stability than Post inversion formula). The Post method approximates $g(\tau^*)$ as:

$$g(\tau^*) \approx \frac{(-1)^k}{k!} \left(\frac{k}{\tau^*}\right)^{k+1} \frac{d^k}{d\sigma^k} V^{\pi}_{\gamma}(s) \bigg|_{\sigma=k/\tau^*},\tag{6}$$

126 where k is a parameter chosen to balance accuracy and computational stability as discussed in the 127 next section, and $\gamma = e^{-\sigma}$.

128 Once $g(\tau^*)$ is estimated, we can compute the value of the policy under any arbitrary temporal 129 discounting function. Examples include finite horizon: for a horizon *T*, the cumulative expected 130 reward is:

$$V_{0\to T}^{\pi}(s) = \sum_{\tau^*=0}^{T} g(\tau^*), \tag{7}$$

131 and hyperbolic and power-law discounting: with a discount function $f(\tau^*) = \frac{1}{(1+K\tau^*)^{\beta}}$, the value 132 is:

$$V_{\text{hyperbolic}}^{\pi}(s) = \sum_{\tau^*=0}^{\infty} \frac{1}{(1+K\tau^*)^{\beta}} g(\tau^*), \tag{8}$$

133 where K > 0 controls the discounting rate.

Fig.1a shows how value increases as a function of τ^* : without temporal discounting the value reaches the true magnitude of the reward. Fig. 2 illustrate the impact of finite horizon on the value for different distances from the goal.



Figure 1: Value function as a result of integrating expected reward at future time steps for different values of reward (with k = 10) (a) and parameter k (b). The reward is located at a distance of 10 steps from the agent.



Figure 2: Value of $\int_{\tau^*=0}^{\tau^*=m} g(\tau^*)$ across varying deadlines (m) and distance from rewarding goal state (τ) for a terminal reward of magnitude 1 for k = 30.

137 Overall, this approach eliminates the need to retrain the policy for different temporal preferences,

138 offering a flexible framework for post-hoc adaptation. We next discuss the properties of the discrete

139 approximation and implication of the log-compression of the timeline.

140 2.2 Emergence of Log-Compression in the Estimation of Future Time

141 In this section, we investigate how the discrete approximation employed to estimate the expected 142 reward function $g(\tau^*)$ naturally yields a log-compressed representation of future time. This property 143 arises from the mathematical structure of the inverse transform approximation and carries important 144 implications for an agent's temporal perception and decision-making.

To illustrate the emergence of log-compression we take advantage of linearity of the Laplace and inverse Laplace transform and first derive its impulse response and later generalize that to arbitrary reward functions. To compute the impulse response, we consider an environment where a single reward of magnitude 1 is positioned at a temporal distance τ from the current state, indicating the number of steps to reach the reward. The expected reward at a future time step τ^* is:

$$g(\tau, \tau^*) = \frac{1}{\tau} \frac{k^{k+1}}{k!} \left(\frac{\tau}{\tau^*}\right)^{k+1} e^{-k\left(\frac{\tau}{\tau^*}\right)}.$$
(9)

The result has a form of gamma distribution (Weisstein, 2004), and has been widely used in learning
temporal relationships (Hopfield & Tank, 1990; Grossberg & Schmajuk, 1989; Shankar & Howard,
2012; Jacques et al., 2021).

As shown by Shankar & Howard (2012) taking the partial derivative of the impulse response of $g(\tau, \tau^*)$ with respect to τ^* and setting it to zero, we find that it is a unimodal function of time t with a maximum at $\tau^* \frac{k}{k+1}$ (in the limit when $k \to \infty$, the peak is exactly at τ^* , see Appendix S3 for derivation). Expected reward as a function of future time for reward magnitude of 1 and distance of 25 steps is shown in Fig. 3 for different values of k. We used 50 log-spaced τ^* values ranging from 1 to 50.

To demonstrate that the representation is indeed log-compressed, we express the width of the unimodal impulse response of $g(\tau, \tau^*)$ through the coefficient of variation $CV = \frac{1}{\sqrt{k+1}}$, which is a ratio of standard deviation and mean (see Appendix S4 for derivation). Critically, CV does not



Figure 3: Expected reward at future time steps for different values of k. The reward of magnitude 1 is presented at distance 25. Increasing k results in more narrow peaks that approach closer to the value of the reward (since $\tau_{\text{peak}}^* = \tau k/(k+1)$).

depend on t and τ^* , implying that the width of the unimodal basis functions increases linearly with their peak time indicating log-compression.

164 If discrete time steps $\tau_1^*, \tau_2^*, \ldots, \tau_m^*$ are log-spaced, i.e., $\tau_i^* = e^{ci}$ (where *c* sets the spacing), then 165 $\log(\tau_i^*) = ci$, yielding constant intervals on a logarithmic scale. Peaks at $\tau_{\text{peak}}^* = \frac{k}{k+1}\tau$ are also 166 log-spaced, as $\log(\tau_{\text{peak}}^*) = \log\left(\frac{k}{k+1}\right) + \log(\tau)$. With constant CV, peak widths scale with τ_{peak}^* , 167 maintaining consistent relative width on a log scale, thus compressing distant time steps relative to 168 near ones.

For temporal discounting with a finite horizon we provide lemmas and proofs for analytical solutions for three different integral bounds: 0 to m (Fig. 4a), m to n (Fig. 4b) and m to ∞ (Fig. 4c). Proofs for the lemmas are in Supplemental Information (Sec. S2).

172 **Lemma 1:** The value of a state at distance τ from the reward given the deadline *m* can be computed 173 as:

$$\int_{0}^{m} g(\tau, \tau^{*}) d\tau^{*} = e^{-k\frac{(\tau)}{m}} \sum_{i=0}^{k-1} \frac{(k\frac{(\tau)}{m})^{i}}{i!}.$$
(10)

174 **Lemma 2:** The state-value perceived by the agent between steps $\tau^* = m$ and $\tau^* = n$ in the future 175 where m < n when the agent is at distance τ from the reward can be computed as follows:

$$S(n,c,x) = e^{cx} \sum_{i=0}^{n} (-1)^{n-i} \frac{n!}{i!c^{n-i+1}} x^{i}$$
$$\int_{m}^{n} g(\tau,\tau^{*}) d\tau^{*} = (-1)^{k+1} (\frac{k^{k+1}}{k!}) [S(k-1,k,\frac{-\tau}{n}) - S(k-1,k,\frac{-\tau}{m})]$$
(11)

176 **Lemma 3:** The value perceived by the agent between step m and ∞ in the future when the agent is 177 at a distance τ from the reward can be computed as follows:

$$\int_{m}^{\infty} g(\tau, \tau^{*}) d\tau^{*} = 1 - e^{-k \frac{(\tau)}{m}} \sum_{i=0}^{k-1} \frac{(k \frac{(\tau)}{m})^{i}}{i!}$$
(12)



Figure 4: Illustration of expected reward as a function of the future time integrated over different bounds: 0 to m (a), m to n (b) and m to ∞ (c). In this example, the reward of magnitude 1 is located 25 steps from the current state (marked as 0). The total area under each curve corresponds to the magnitude of the reward. Analytical derivation for each integral is provided in the Supplemental Information (Sec. S2).

178 2.3 Policy Selection Under Arbitrary Discounting Functions

179 In this subsection, we extend our framework to select the optimal policy from a set of precomputed 180 policies under any user-specified temporal discounting function. This method leverages the expected 181 reward function $g(\tau^*)$ introduced earlier, enabling adaptability to diverse temporal preferences with-182 out requiring retraining or additional environment interactions.

183 Given a set of precomputed policies $\Pi = \{\pi_1, \pi_2, \dots, \pi_p\}$, we assume each policy π_k has been 184 evaluated to compute state-value functions $V_{\gamma_i}^{\pi_k}(s)$ for a discrete set of discount factors Γ . As 185 described above, we approximate the expected reward function $g_k(\tau^*) = \mathbb{E}_{\pi_k}[r_{t+\tau^*} \mid s_t = s]$ for 186 each policy π_k using the inverse transform applied to the sampled values $V_{\gamma_i}^{\pi_k}(s)$.

For a user-specified temporal discounting function $f : \mathbb{N} \to \mathbb{R}^+$, which assigns weights to rewards based on their temporal distance τ^* , we compute the value of each policy π_k as:

$$V_f^{\pi_k}(s) = \sum_{\tau^*=0}^{\infty} f(\tau^*) g_k(\tau^*).$$
(13)

189 This value represents the expected cumulative reward under policy π_k adjusted by the discounting 190 function *f*. The optimal policy from the set Π is then selected as:

$$\pi^* = \arg\max_{\pi_k \in \Pi} V_f^{\pi_k}(s). \tag{14}$$

191 This approach allows the agent to adapt to arbitrary discounting schemes post-training. Under ide-192 alized conditions—where the set of discount factors Γ is continuous and $g_k(\tau^*)$ is perfectly re-193 covered—the selected policy π^* is optimal within Π (not globally optimal) for the given f (see 194 Appendix S5 for proof of optimality).

195 In practice, Γ is finite, and $g_k(\tau^*)$ is approximated as $\hat{g}_k(\tau^*)$. Thus, we compute $\hat{V}_f^{\pi_k}(s) = \sum_{\tau^*=0}^{\infty} f(\tau^*) \hat{g}_k(\tau^*)$, and the selected policy is optimal with respect to these approximations.

197 2.4 Dynamic Policy Selection

198 In this subsection, we present an extension to our framework that enables *dynamic policy selection* 199 at each time step. This approach allows the agent to adapt its strategy based on the current state and 200 a user-specified temporal discounting function f. At each time step t, the agent, in state s_t , computes the value of each policy π_k from that state under the discounting function f. This value is given by:

$$V_f^{\pi_k}(s_t) = \sum_{\tau^*=0}^{\infty} f(\tau^*) g_k(\tau^* \mid s_t),$$
(15)

where $g_k(\tau^* \mid s_t)$ represents the expected reward at time $t + \tau^*$ under policy π_k , starting from s_t . The agent selects the policy π_{k^*} that maximizes this value:

$$\pi_{k^*} = \arg\max_{\pi_k \in \Pi} V_f^{\pi_k}(s_t).$$
(16)

The action at s_t is then $a_t = \pi_{k^*}(s_t)$. Upon transitioning to the next state s_{t+1} , the agent repeats this process, recomputing $V_f^{\pi_k}(s_{t+1})$ for all $\pi_k \in \Pi$ and selecting the optimal policy for s_{t+1} , which may differ from the previous choice.

This dynamic process constructs a hybrid policy π_{hybrid} , where the action at each state s_t is determined by the policy π_{k^*} that maximizes $V_f^{\pi_k}(s_t)$ at that time. The hybrid policy π_{hybrid} is not necessarily an element of Π , as it may combine actions from multiple policies depending on the state, where $\pi_{k^*} = \arg \max_{\pi_k \in \Pi} V_f^{\pi_k}(s_t)$. By greedily selecting the action from the policy that maximizes the value at each state, π_{hybrid} locally optimizes the expected cumulative reward under f. In environments where different policies perform better in different regions of the state space, this approach can outperform any single policy in Π .

215 We illustrate this approach through a grid-world example (Fig. 5). The agent starts from a corner of the grid. Five rewards are placed at different locations in the environment such that rewards with a 216 217 larger magnitude are placed further from the agent's start location. We first compute five different 218 policies and corresponding values such that each of the policies leads the agent to a different reward. 219 When the agent starts navigating, it is given a number of steps to complete the task (finite horizon). 220 The agent follows Alg. 1 and at each state computes the expected future reward for each of 5 policies 221 for the available remaining number of steps. The agent always selects a policy that leads to the 222 largest reward and takes the corresponding action. In our example, the agent always reached the 223 largest possible reward given the available number of steps (note that this is not guaranteed, as we 224 investigate in the following subsections). We emphasize that agents based on traditional exponential 225 discounting with a single value of γ would always converge to the same reward since those agents 226 do not have the capability to take into account the available horizon.



Figure 5: Paths chosen by the model for a varying number of available steps, T. The colors represent reward locations and correspond to rewards of different magnitudes. Higher rewards are located further. Given a longer horizon, the agent chooses rewards located further to obtain a larger reward.

227 2.5 Efficient Performance in Scale-Free Environments

In environments lacking a single dominant timescale—often called *scale-free*—reward timing may follow a power-law pattern in its temporal correlations, for instance $\mathbb{E}[r_t r_{t+\tau^*}] \propto (\tau^*)^{-\alpha}$ with

Algorithm 1 Dynamic Policy Selection

1: Input: Set of policies $\Pi = \{\pi_1, \pi_2, \dots, \pi_p\}$, temporal discounting function f, initial state s_0 , maximum time steps T, horizon for each time step h_0, h_1, \dots, h_{T-1} , approximated expected reward functions $g_k(\tau^* \mid s)$ for each $\pi_k \in \Pi$ and state s 2: Set $s \leftarrow s_0$ 3: Set $t \leftarrow 0$ 4: while t < T and s is not terminal do for each $\pi_k \in \Pi$ do 5: Compute $V_f^{\pi_k}(s) = \sum_{\tau^*=0}^{h_t} f(\tau^*) g_k(\tau^* \mid s)$ 6: 7: end for Select $\pi_{k^*} = \arg \max_{\pi_k \in \Pi} V_f^{\pi_k}(s)$ 8: Take action $a = \pi_{k^*}(s)$ 9: Observe next state s' and reward r10: Set $s \leftarrow s'$ 11: Set $t \leftarrow t+1$ 12: 13: end while

230 $\alpha > 0$. This implies that reward timing variances scale with the square of their mean, $Var(\tau) \propto$ 231 τ^2 —a relationship observed in interval-timing tasks in humans and other animals.

To illustrate this, suppose the agent can pursue two rewards: one at $\tau = 10$ steps (with variance ± 1 step), and another at $\tau = 100$ steps (with variance ± 10 steps). In our log-compressed timeline, each reward is represented by a unimodal function centered near its mean arrival time, whose "width" is proportional to that mean. Thus, the reward at 10 steps exhibits a narrow peak, while the reward at 100 steps has a broader peak. When integrating the expected reward up to a finite horizon T = 50, that is

$$V_{0\to T}^{\pi}(s) = \sum_{\tau^*=0}^{T} g(\tau^*), \qquad (17)$$

the entire distribution for the near reward lies within the 50-step horizon. In contrast, only a portion of the broader distribution at 100 steps overlaps the first 50 steps. However, if that distant reward is large enough, the probability of it arriving earlier than 100 steps can contribute a higher total expected value than the smaller near reward. By preserving the full variance structure around each expected time, the agent's scale-adaptive representation captures this partial overlap automatically.

The *efficiency* of such a log-compressed representation can be framed in terms of minimizing the maximum *relative* error under a scale-invariant prior $p(\tau) \propto 1/\tau$. Specifically, for the expected squared relative error,

$$\operatorname{Error} = \mathbb{E}\left[\left(\frac{\hat{g}(\tau^*) - g(\tau^*)}{g(\tau^*)}\right)^2\right],\tag{18}$$

246 this approach ensures that

$$\frac{\sqrt{\operatorname{Var}(\hat{g}(\tau^*))}}{g(\tau^*)} = \operatorname{CV} = \frac{1}{\sqrt{k+1}},$$
(19)

remains constant across τ^* . By contrast, a linear-time representation with uniform variance allocates precision in a way that either under-resolves near-future rewards or over-resolves distant ones. The key is that the logarithmic transform *stabilizes* variance, because $Var(log(\tau)) \approx (\frac{1}{\tau})^2 Var(\tau)$ is nearly constant when $Var(\tau) \propto \tau^2$. Consequently, the log compression allocates resources proportionally to the relevant timescale, ensuring uniform relative precision for near and distant rewards—particularly advantageous in scale-free settings and under limited time horizons.

253 2.6 Controlling the Precision of Reward Estimation

The parameter k in the log-compressed representation governs the precision of reward estimation, with the coefficient of variation given by CV. Larger values of k reduce the CV (Fig. 3), enhancing 256 accuracy but increasing computational demands and impact stability of the inverse transform (since 257 k controls the order of the numerical derivative in the Post inversion, Eq. 6). The change in expected 258 future reward as a function of future time for different values of k is shown in Fig. 1b. The choice 259 of k can impact the choice that the agent makes. This is illustrated in Fig 6a where with k = 3 agent prefers a closer but smaller reward and in Fig 6b where the only change is that k = 80 but the agent 260 261 now prefer more distant but larger reward. This flexibility in k reflects a tunable trade-off between 262 precision and computational cost, allowing the representation to adapt to different environments or 263 resource constraints. We emphasize that this does not undermine the representation's optimality, 264 which ensures uniform relative precision across τ^* for any chosen k. This variability in k also offers 265 a parallel to human timing behavior, where the scalar property implies a constant CV that differs 266 across individuals.



Figure 6: Comparison of value function for different values of k in the presence of rewards of different magnitudes at different distances.

267 **3 Discussion**

268 This paper introduces a framework for reinforcement learning (RL) that separates temporal discount-269 ing from the training phase, allowing adjustments to the discount function after training without the 270 need for retraining. By employing a log-compressed representation of expected future rewards, de-271 rived from value functions computed across a range of discount factors, the approach enables agents 272 to adapt flexibly to various temporal preferences—such as exponential, hyperbolic, or finite-horizon 273 discounting. Additionally, we propose a dynamic policy selection mechanism, utilizing a precom-274 puted library of policies to construct a hybrid strategy that adjusts to shifting objectives at runtime. 275 These contributions enhance the adaptability of RL systems, which could prove useful in applica-276 tions where temporal preferences vary, such as robotics or financial decision-making. Furthermore, 277 the log-compressed representation aligns conceptually with human perception of time, as described 278 by the Weber-Fechner law, suggesting a potential avenue for modeling human-like decision pro-279 cesses in artificial agents.

280 Despite its potential, the framework has several limitations that warrant consideration:

Approximation Errors: The reliance on an approximated inverse Laplace transform, based on a finite set of discount factors, introduces possible inaccuracies in reconstructing the desired discount function. The quality of this approximation depends heavily on the number and distribution of discount factors chosen, which may not generalize across all environments.

Computational Overhead: Reconstructing the reward timeline and selecting policies dynamically at each step can demand significant computational resources, particularly in large or complex state spaces. The feasibility of scaling this approach to high-dimensional tasks, such as continuous control, has yet to be demonstrated.

- Policy Library Dependence: The success of dynamic policy selection hinges on the diversity and quality of the precomputed policy library. If the library lacks policies suited to a specific discount
- 291 objective, performance may fall short compared to a policy trained specifically for that purpose.
- This work offers a step toward greater flexibility in RL, enabling agents to adapt to varying temporal preferences and objectives in a manner that may resonate with human perception.

294 **References**

- George Ainslie. Specious reward: a behavioral theory of impulsiveness and impulse control. *Psy- chological Bulletin*, 82(4):463, 1975.
- Rui Cao, Ian M Bright, and Marc W Howard. Ramping cells in the rodent medial prefrontal cortex
 encode time to past and future events via real laplace transform. *Proceedings of the National Academy of Sciences*, 121(38):e2404169121, 2024.
- 300 Gustav Fechner. *Elements of Psychophysics. Vol. I.* Houghton Mifflin, 1860/1912.
- William Fedus, Carles Gelada, Yoshua Bengio, Marc G Bellemare, and Hugo Larochelle. Hyper bolic discounting and learning over multiple horizons. *arXiv preprint arXiv:1902.06865*, 2019.
- Shane Frederick, George Loewenstein, and Ted O'donoghue. Time discounting and time preference:
 A critical review. *Journal of economic literature*, 40(2):351–401, 2002.
- John Gibbon. Scalar expectancy theory and weber's law in animal timing. *Psychological review*, 84 (3):279, 1977.
- Leonard Green and Joel Myerson. Exponential versus hyperbolic discounting of delayed outcomes:
 Risk and waiting time. *American Zoologist*, 34(4):496–505, 1994.
- Stephen Grossberg and Nestor A Schmajuk. Neural dynamics of adaptive timing and temporal
 discrimination during associative learning. *Neural networks*, 2(2):79–102, 1989.
- John J Hopfield and David W Tank. Neural computation by time concentration, June 26 1990. US
 Patent 4,937,872.
- Marc W Howard and Karthik H Shankar. Neural scaling laws for an uncertain world. *Psychological review*, 125(1):47, 2018.
- Brandon Jacques, Zoran Tiganj, Marc W Howard, and Per B Sederberg. Deepsith: Efficient learn ing via decomposition of what and when across time scales. *Advances in neural information processing system*, 2021.
- Zeb Kurth-Nelson and A David Redish. Temporal-difference reinforcement learning with distributed
 representations. *PLoS One*, 4(10):e7362, 2009.
- David Laibson. Golden eggs and hyperbolic discounting. *The Quarterly Journal of Economics*, 112 (2):443–477, 1997.
- Paul Masset, Pablo Tano, HyungGoo R Kim, Athar N Malik, Alexandre Pouget, and Naoshige
 Uchida. Multi-timescale reinforcement learning in the brain. *bioRxiv*, 2023.
- András Mészáros and Miklós Telek. Concentrated matrix exponential distributions with real eigenvalues. *Probability in the Engineering and Informational Sciences*, 36(4):1171–1187, 2022.
- Ida Momennejad and Marc W Howard. Predicting the future with multi-scale successor representations. *BioRxiv*, pp. 449470, 2018.
- Emil L Post. Generalized differentiation. *Transactions of the American Mathematical Society*, 32
 (4):723–781, 1930.

- A David Redish. Addiction as a computational process gone awry. Science, 306(5703):1944–1947, 330 2004. 331
- Matthias Schultheis, Constantin A. Rothkopf, and Heinz Koeppl. Reinforcement learning with non-332 exponential discounting. In Advances in Neural Information Processing Systems, 2022. 333
- Karthik H Shankar and Marc W Howard. A scale-invariant internal representation of time. Neural 334 Computation, 24(1):134-193, 2012. 335
- 336 Karthik H Shankar and Marc W Howard. Optimally fuzzy temporal memory. The Journal of 337 Machine Learning Research, 14(1):3785-3812, 2013.
- 338 Richard S. Sutton and Andrew G. Barto. Reinforcement Learning: An Introduction. MIT Press, 339 2018.
- 340 Yunhao Tang, Mark Rowland, Rémi Munos, and Michal Valko. Taylor expansion of discount factors. 341 In International Conference on Machine Learning, pp. 10130–10140. PMLR, 2021.
- 342 Pablo Tano, Peter Dayan, and Alexandre Pouget. A local temporal difference code for distributional 343 reinforcement learning. Advances in neural information processing systems, 33:13662–13673, 344 2020.
- 345 Zoran Tiganj, Samuel J Gershman, Per B Sederberg, and Marc W Howard. Estimating scaleinvariant future in continuous time. Neural Computation, 31(4):681-709, 2019. 346
- 347 Zoran Tiganj, Inder Singh, Zahra G Esfahani, and Marc W Howard. Scanning a compressed ordered representation of the future. Journal of Experimental Psychology: General, pp. In Press, 2022. 348
- 349 Eric W Weisstein. Gamma distribution. https://mathworld. wolfram. com/, 2004.
- Jason T Wilkes. Reverse first principles: Weber's law and optimality in different senses. PhD thesis, 350
- 351 UNIVERSITY OF CALIFORNIA, SANTA BARBARA, 2015.

352 353

354

Supplementary Materials

The following content was not necessarily subject to peer review.

355 S1 Numerical Approximations of the Inverse Laplace Transform

In this section, we describe a numerical approximation for computing the inverse Laplace transform
 using the CME-R (Complex Matrix Exponentials with Real eigenvalues) method (introduced by
 Mészáros & Telek (2022)). This method provides better numerical stability when compared to Post
 method (Post, 1930) (Fig S1).

360 Let $\beta, \eta \in \mathbb{R}^M$, where $\tau_i^* \in \mathbb{R}$, $\sigma_{ij} = \frac{\beta_j}{\tau_i^*}$ and $\gamma_{ij} = e^{-\sigma_{ij}}$. The approximation takes the form:

$$g(\tau_i^*) \approx \sum_{k=1}^M \frac{\eta_k}{\tau_i^*} V_{\gamma_{i,k}}^{\pi}$$
(20)

Here, M represents the maximum number of function evaluations and serves a similar role to the parameter k in the post-inverse method. Both η and β are real-valued tensors. It is worth noting that Equation 20 requires higher numerical precision as M increases to obtain accurate inverse Laplace transform results. The detailed methodology for computing the parameters η and β is thoroughly

365 documented in the work of Mészáros & Telek (2022).



Figure S1: Comparing expected return between CME-R inverse, Offline Post Inverse and Online Post Inverse. While this figure represents a general property of the inverse Laplace transform for a time-series input, in our case, the input (blue) is a time-varying reward function. Despite high temporal variability of the input (reward), CME-R (orange) is able to construct a fateful log-compressed estimate with temporal resolution gradually decaying from 0 (the current state) towards the future. Offline Post inverse (red) is an analytically computed offline solution that convolves the input function with the set of impulse responses defined in Eq. 9. (It serves as a ground truth since it does not suffer from numerical issues.) A high match between CME-R and the Offline Post inverse (computed online using Eq. 6) suffers from numerical instabilities when the input signal has high temporal variability (green).

366 S2 Lemmas and Proofs

We define $\tau^*, \tau \in [0, \infty]$. τ is the distance from the reward and τ^* refers to the number of steps in the future from the agent's current position and $\sigma = \frac{k}{\tau^*}$. If we assume there is a single terminal 369 reward with magnitude 1 in an environment, then the value of state s at distance τ from the reward 370 is given by $\gamma^{\tau} = (e^{-\sigma})^{\tau} = e^{-\sigma\tau}$. Then we can compute the expected reward at time step τ_i^* in the 371 future as follows:

$$g(\tau, \tau^*) = \mathcal{L}^{-1} \{ V^{\pi}_{\gamma=e^{-\sigma\tau}}(s) \}$$

= $\mathcal{L}^{-1} \{ e^{-\sigma\tau} \}$
= $\frac{(-1)^k}{k!} \sigma^{k+1} \frac{d^k}{d\tau^k} e^{-\sigma\tau}$
= $\frac{(-1)^k}{k!} \sigma^{k+1} (-\tau)^k e^{-\sigma\tau}$
= $\frac{1}{\tau} \frac{k^{k+1}}{k!} (\frac{\tau}{\tau^*})^{k+1} e^{-k(\frac{\tau}{\tau^*})}$

372 **Lemma 1:** The value of a state at distance τ from the reward given the deadline m can be computed 373 as:

$$\int_{0}^{m} g(\tau, \tau^{*}) d\tau^{*} = e^{-k\frac{(\tau)}{m}} \sum_{i=0}^{k-1} \frac{(k\frac{(\tau)}{m})^{i}}{i!}.$$
(21)

374 **Proof:**

$$\begin{split} \int_0^m g(\tau,\tau^*) d\tau^* &= \int_0^m \frac{1}{\tau} \frac{k^{k+1}}{k!} (\frac{\tau}{\tau^*})^{k+1} e^{-k(\frac{\tau}{\tau^*})} d\tau^* \\ &= \frac{1}{\tau} \frac{k^{k+1}}{k!} \int_0^m (\frac{\tau}{\tau^*})^{k+1} e^{-k(\frac{\tau}{\tau^*})} d\tau^* \end{split}$$

Substituting $u = \frac{\tau}{\tau^*}$,

$$= -\frac{k^{k+1}}{k!} \int_{\infty}^{\frac{\tau}{m}} u^{k-1} e^{-ku} du$$

$$= \frac{k^{k+1}}{k!} \int_{\frac{\tau}{m}}^{\infty} u^{k-1} e^{-ku} du$$

$$= \frac{k^{k+1}}{k!} \left[\int_{0}^{\infty} u^{k-1} e^{-ku} du - \int_{0}^{\frac{(\tau)}{m}} u^{k-1} e^{-ku} du \right]$$

Using $\int_0^\infty y^m e^{-ky} dy = \frac{\Gamma(m+1)}{k^{m+1}},$

$$= \frac{k^{k+1}}{k!} \left[\frac{\Gamma(k)}{k^k} - \int_0^{\frac{(\tau)}{m}} u^{k-1} e^{-ku} du \right]$$
$$= 1 - \frac{k^{k+1}}{k!} \int_0^{\frac{(\tau)}{m}} u^{k-1} e^{-ku} du$$

$$\begin{aligned} \text{Using } \int_0^b x^n e^{-ax} dx &= \frac{n!}{a^{n+1}} \left[1 - e^{-ab} \sum_{i=0}^{i=n} \frac{(ab)^i}{i!} \right], \\ &= 1 - \frac{k^{k+1}}{k!} \frac{(k-1)!}{k^k} \left[1 - e^{-k\frac{(\tau)}{m}} \sum_{i=0}^{i=k-1} \frac{(k\frac{(\tau)}{m})^i}{i!} \right] \\ &= 1 - 1 \left[1 - e^{-k\frac{(\tau)}{m}} \sum_{i=0}^{i=k-1} \frac{(k\frac{(\tau)}{m})^i}{i!} \right] \\ &= e^{-k\frac{(\tau)}{m}} \sum_{i=0}^{i=k-1} \frac{(k\frac{(\tau)}{m})^i}{i!} \end{aligned}$$

1375 **Lemma 2:** The state-value perceived by the agent between steps $\tau^* = m$ and $\tau^* = n$ in the future 1376 where m < n when the agent is at distance τ from the reward can be computed as follows:

$$S(n,c,x) = e^{cx} \sum_{i=0}^{n} (-1)^{n-i} \frac{n!}{i!c^{n-i+1}} x^{i}$$
$$\int_{m}^{n} g(\tau,\tau^{*}) d\tau^{*} = (-1)^{k+1} (\frac{k^{k+1}}{k!}) [S(k-1,k,\frac{-\tau}{n}) - S(k-1,k,\frac{-\tau}{m})].$$
(22)

377 **Proof:**

$$\begin{split} \int_{m}^{n} g(\tau,\tau^{*}) d\tau^{*} &= \int_{m}^{n} \frac{1}{\tau} \frac{k^{k+1}}{k!} (\frac{\tau}{\tau^{*}})^{k+1} e^{-k(\frac{\tau}{\tau^{*}})} d\tau^{*} \\ &= (-1)^{k+1} \frac{1}{\tau} \frac{k^{k+1}}{k!} \int_{m}^{n} (\frac{-\tau}{\tau^{*}})^{k+1} e^{-k(\frac{\tau}{\tau^{*}})} d\tau^{*} \end{split}$$

378 Substituting $u = \frac{-\tau}{\tau^*}$,

$$= (-1)^{k+1} \frac{1}{\tau} \frac{k^{k+1}}{k!} \int_{\frac{-\tau}{m}}^{\frac{-\tau}{n}} u^{k+1} e^{ku} \frac{\tau}{u^2} du$$
$$= (-1)^{k+1} \frac{1}{\tau} \frac{k^{k+1}}{k!} \int_{\frac{-\tau}{m}}^{\frac{-\tau}{n}} u^{k-1} e^{ku} du$$

379 Using
$$S(n,c,x) = \int x^n e^{cx} dx = e^{cx} \sum_{i=0}^n (-1)^{n-i} \frac{n!}{i!c^{n-i+1}} x^i,$$

$$= (-1)^{k+1} \left(\frac{k^{k+1}}{k!}\right) \left[S(k-1,k,\frac{-\tau}{n}) - S(k-1,k,\frac{-\tau}{m})\right]$$

Lemma 3: The value perceived by the agent between step m and ∞ in the future when the agent is at a distance τ from the reward can be computed as follows:

$$\int_{m}^{\infty} g(\tau, \tau^{*}) d\tau^{*} = 1 - e^{-k \frac{(\tau)}{m}} \sum_{i=0}^{k-1} \frac{(k \frac{(\tau)}{m})^{i}}{i!}.$$
(23)

382 **Proof:**

$$\begin{split} \int_{m}^{\infty} g(\tau,\tau^{*}) d\tau^{*} &= \int_{0}^{\infty} g(\tau^{*}) d\tau^{*} - \int_{0}^{m} g(\tau^{*}) d\tau^{*} \\ &= (\frac{1}{\tau} \frac{k^{k+1}}{k!} \int_{0}^{\infty} (\frac{\tau}{\tau^{*}})^{k+1} e^{-k(\frac{\tau}{\tau^{*}})} d\tau^{*}) - (\int_{0}^{m} g(\tau^{*}) d\tau^{*}) \end{split}$$

Substituting $u = \frac{\tau}{\tau^*}$,

$$= \left(\frac{-k^{k+1}}{k!} \int_{\infty}^{0} u^{k-1} e^{-ku} du\right) - \left(\int_{0}^{m} g(\tau^{*}) d\tau^{*}\right)$$
$$= \left(\frac{k^{k+1}}{k!} \int_{0}^{\infty} u^{k-1} e^{-ku} du\right) - \left(\int_{0}^{m} g(\tau^{*}) d\tau^{*}\right)$$

Using $\int_0^\infty y^m e^{-ky} dy = \frac{\Gamma(m+1)}{k^{m+1}}$,

$$=1-\int_0^m g(\tau^*)d\tau^*$$

Using Lemma 1,

$$= 1 - e^{-k\frac{(\tau)}{m}} \sum_{i=0}^{k-1} \frac{(k\frac{(\tau)}{m})^i}{i!}$$

383 S3 Derivation of the Peak Time of the Impulse Response

$$\frac{\partial \mathcal{T}(\tau,\tau^*)}{\partial \tau^*} = 0.$$

Befine $u = \frac{\tau}{\tau^*}$, so the function becomes $\mathcal{T}(\tau, \tau^*) = \frac{1}{\tau} \frac{k^{k+1}}{k!} u^{k+1} e^{-ku}$. Since τ is fixed, we differentiate with respect to u, where $\tau^* = \frac{\tau}{u}$ and $\frac{d\tau^*}{du} = -\frac{\tau}{u^2}$. Using the chain rule:

$$\frac{\partial \mathcal{T}}{\partial \tau^*} = \frac{\partial \mathcal{T}}{\partial u} \cdot \frac{du}{d\tau^*},$$

386 with $\frac{du}{d\tau^*} = -\frac{\tau}{(\tau^*)^2}$. Compute the derivative:

$$\frac{\partial \mathcal{T}}{\partial u} = \frac{1}{\tau} \frac{k^{k+1}}{k!} \left[(k+1)u^k e^{-ku} - ku^{k+1} e^{-ku} \right] = \frac{1}{\tau} \frac{k^{k+1}}{k!} e^{-ku} u^k \left[(k+1) - ku \right].$$

387 Setting this to zero:

$$(k+1) - ku = 0 \quad \Rightarrow \quad u = \frac{k+1}{k}$$

388 Since $u = \frac{\tau}{\tau^*}$, we solve for τ^* :

$$\tau^* = \frac{k}{k+1}\tau$$

Thus, the peak occurs at $\tau_{\text{peak}}^* = \frac{k}{k+1}\tau$, showing a linear relationship with τ , where the factor $\frac{k}{k+1}$ approaches 1 as k grows, enhancing accuracy.

391 S4 Derivation of CV of the Impulse Response

392 The mean of the impulse response of $g(\tau, \tau^*)$ is:

$$\begin{split} \mu &= \int_0^\infty t \tilde{f}(s;t) dt \\ &= \int_0^\infty t \frac{1}{t} \frac{k^{k+1}}{k!} \left(\frac{t}{\tau^*}\right)^{k+1} e^{-k\frac{t}{\tau^*}} dt \\ &= \frac{k^{k+1}}{k!} \int_0^\infty \left(\frac{t}{\tau^*}\right)^{k+1} e^{-k\frac{t}{\tau^*}} dt \\ &= \tau \frac{k+1}{k}. \end{split}$$

393 The standard deviation of the impulse response of $g(\tau, \tau^*)$ is:

$$\begin{split} \sigma &= \sqrt{\int_0^\infty (t-\mu)^2 \tilde{f}(s;t) dt} \\ &= \sqrt{\int_0^\infty (t-\mu)^2 \frac{1}{t} \frac{k^{k+1}}{k!} \left(\frac{t}{\tau^*}\right)^{k+1} e^{-k\frac{t}{\tau^*}} dt} \\ &= \tau \frac{\sqrt{k+1}}{k}. \end{split}$$

394 Finally, the coefficient of variation is then:

$$CV = \frac{\sigma}{\mu} = \frac{1}{\sqrt{k+1}}.$$

395 S5 Proof of Optimality for Policy Selection

Solution 2013 Consider a finite set of policies Π and a discounting function $f(\tau^*) \ge 0$ with $\sum_{\tau^*=0}^{\infty} f(\tau^*) < \infty$. For each policy $\pi_k \in \Pi$, the true value under f is:

$$V_f^{\pi_k}(s) = \sum_{\tau^*=0}^{\infty} f(\tau^*) g_k(\tau^*),$$

where $g_k(\tau^*) = \mathbb{E}_{\pi_k}[r_{t+\tau^*} | s_t = s]$ is the exact expected reward at time τ^* , fully determined by π_k , \mathcal{P} , and \mathcal{R} . In the limit of continuous Γ , the inverse Laplace transform recovers $g_k(\tau^*)$ exactly from $V_{\gamma^k}(s) = \sum_{\tau^*=0}^{\infty} \gamma^{\tau^*} g_k(\tau^*)$, as $V_{\gamma^k}(s)$ is the Laplace transform of $g_k(\tau^*)$ evaluated at $\sigma = -\ln(\gamma)$. Thus, $V_f^{\pi_k}(s)$ precisely captures the expected discounted reward under discount function f.

403 Since Π is finite, there exists a policy $\pi^* \in \Pi$ such that:

$$V_f^{\pi^*}(s) = \max_{\pi_k \in \Pi} V_f^{\pi_k}(s).$$

404 The selection $\pi^* = \arg \max_{\pi_k \in \Pi} V_f^{\pi_k}(s)$ identifies this policy, proving that π^* is optimal within Π 405 for the given f under these idealized conditions.