

A Additional Experiments

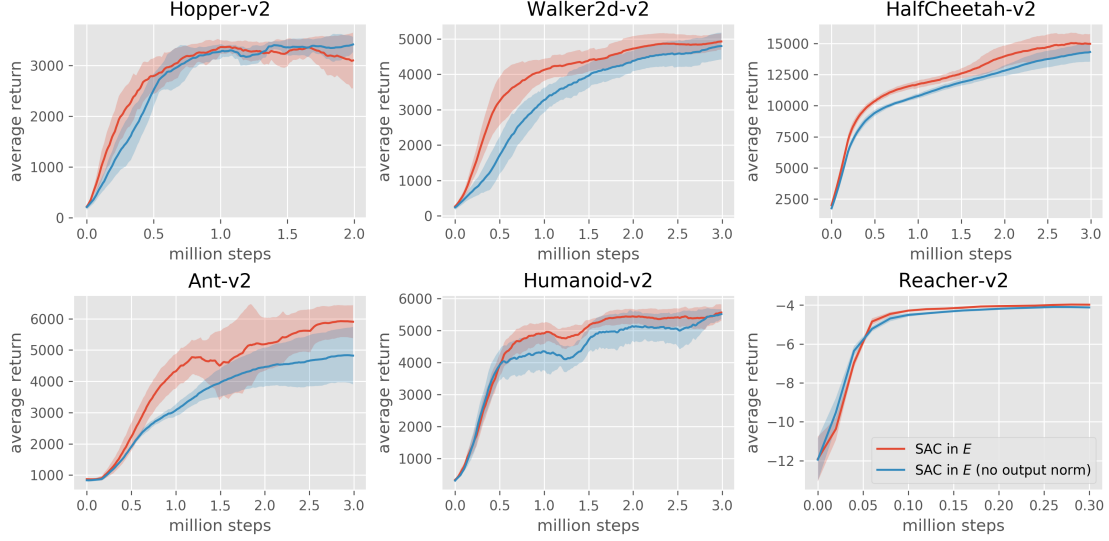


Figure 9: SAC with and without output normalization. SAC in E (no output norm) corresponds to the canonical version presented in Haarnoja et al. (2018a). Mean and 95% confidence interval are computed over eight training runs per environment.

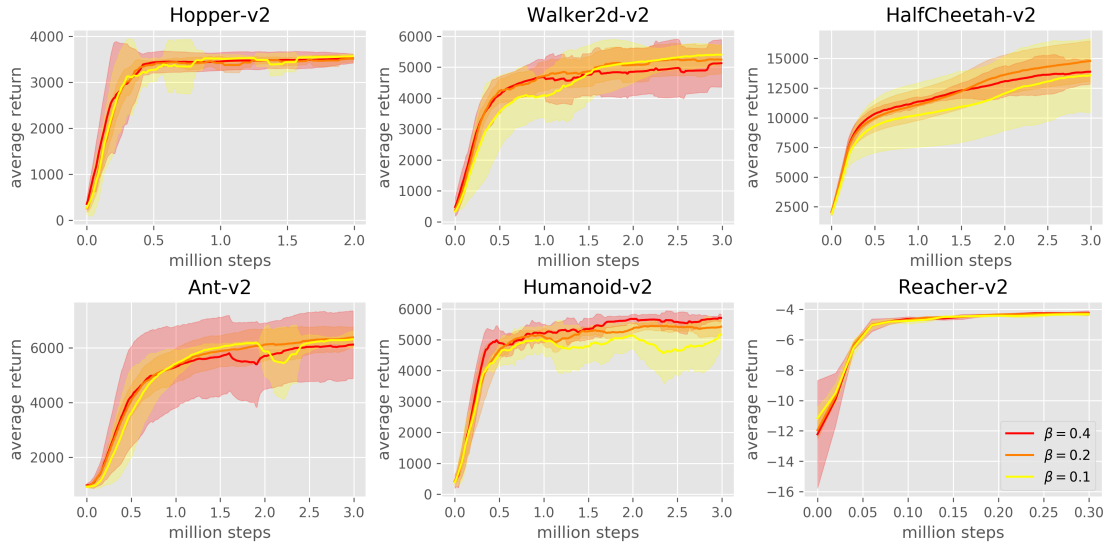


Figure 10: Comparison between different actor loss scales (β). Mean and 95% confidence interval are computed over four training runs per environment.

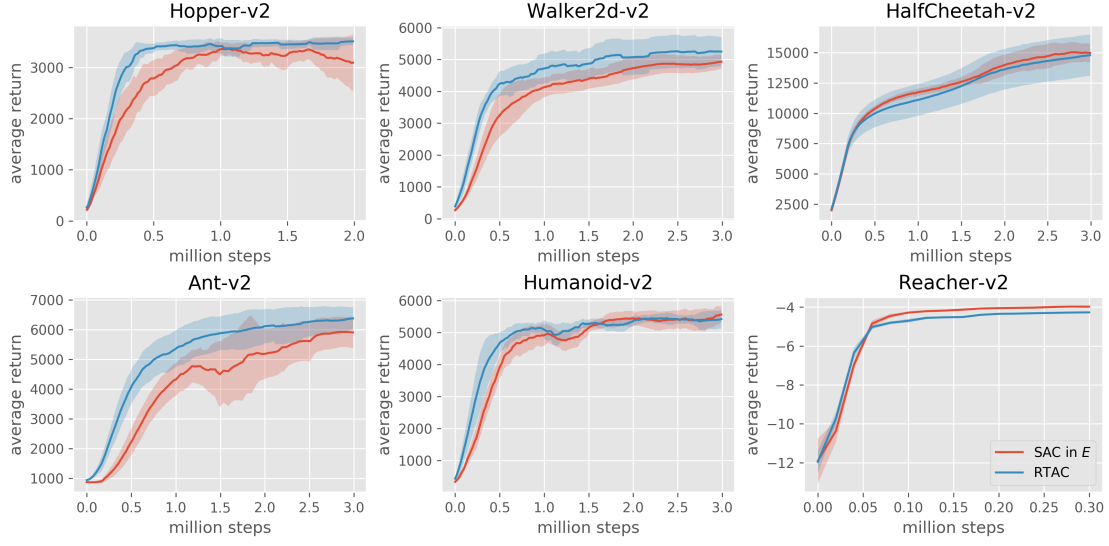


Figure 11: Comparison between RTAC (real-time) and SAC in E (turn-based). Mean and 95% confidence interval are computed over eight training runs per environment.

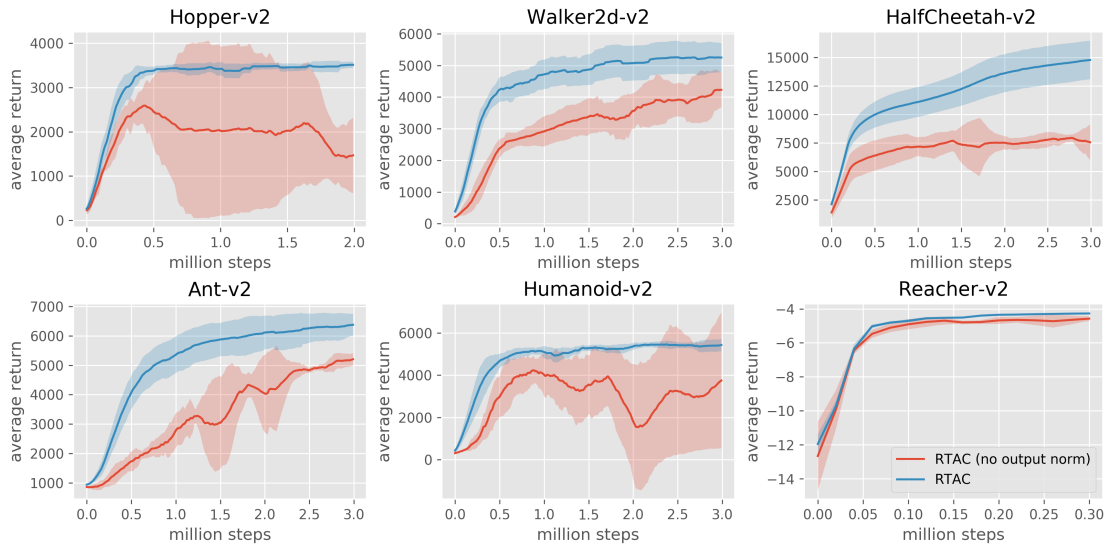


Figure 12: RTAC with and without output normalization. Mean and 95% confidence interval are computed over eight and four training runs per environment, respectively.

B Hyperparameters

Table 1: Hyperparameters

Name	RTAC	SAC
optimizer	Adam	Adam (Kingma & Ba, 2014)
learning rate	0.0003	0.0003
discount (γ)	0.99	0.99
hidden layers	2	2
units per layer	256	256
samples per minibatch	256	256
target smoothing coefficient (τ)	0.005	0.005
gradient steps / environment steps	1	1
reward scale	5	5
entropy scale (α)	1	1
actor-critic loss factor (β)	0.2	-
Pop-Art alpha	0.0003	-
start training after	10000	10000 steps

C Proofs

Theorem 1. ⁴ A policy $\pi : A \times \mathbf{X} \rightarrow \mathbb{R}$ interacting with $RTMDP(E)$ in the conventional, turn-based manner gives rise to the same Markov Reward Process as π interacting with E in real-time, i.e.

$$RTMRP(E, \pi) = TBMRP(RTMDP(E), \pi). \quad (3)$$

Proof. For any environment $E = (S, A, \mu, p, r)$, we want to show that the two above MRPs are the same. Per Def. 2 and 4 for $TBMRP(RTMDP(E), \pi)$ we have

$$\begin{aligned}
(1) \text{ state space} & S \times A, \\
(2) \text{ initial distribution} & \mu(s)\delta(a - c), \\
(3) \text{ transition kernel} & \int_A p(s_{t+1}|s_t, a_t)\delta(a_{t+1} - \mathbf{a}) \pi(\mathbf{a}|s_t, a_t) d\mathbf{a}, \\
(4) \text{ state-reward function} & \int_A r(s, a) \pi(\mathbf{a}|s_t, a_t) d\mathbf{a}.
\end{aligned}$$

The transition kernel, using the definition of the Dirac delta function δ , can be simplified to

$$p(s_{t+1}|s_t, a_t) \int_A \delta(a_{t+1} - \mathbf{a}) \pi(\mathbf{a}|s_t, a_t) d\mathbf{a} = p(s_{t+1}|s_t, a_t) \pi(a_{t+1}|s_t, a_t). \quad (15)$$

The state-reward function can be simplified to

$$r(s_t, a_t) \int_A \pi(\mathbf{a}|\mathbf{x}) d\mathbf{a} = r(s_t, a_t). \quad (16)$$

It should now be easy to see how the elements above match $RTMRP(E, \pi)$, Def. 3. \square

⁴All proofs are in Appendix C.

Theorem 2. A policy $\pi(\mathbf{a}|s, b, a) = \pi(\mathbf{a}|s)$ interacting with $TBMDP(E)$ in real time, gives rise to a Markov Reward Process that contains (Def. 10) the MRP resulting from π interacting with E in the conventional, turn-based manner, i.e.

$$TBMRP(E, \pi) \propto RTMRP(TBMDP(E), \pi) \quad (4)$$

Proof. Given MDP $E = (S, A, \mu, p, r)$, we have $\Psi = (Z, \nu, \sigma, \bar{\rho}) = RTMRP(TBMDP(E), \pi)$ with

$$(1) \text{ state space} \quad Z = S \times \{0, 1\} \times A, \quad (17)$$

$$(2) \text{ initial distribution} \quad \nu(s, b, a) = \mu(s) \delta(b) \delta(a - c), \quad (18)$$

$$(3) \text{ transition kernel} \quad \sigma(s_{t+1}, b_{t+1}, a_{t+1} | s_t, b_t, a_t) \quad (19)$$

$$= \begin{cases} \delta(s_{t+1} - s_t) \delta(b_{t+1} - 1) \pi(a_{t+1} | s_t) & \text{if } b_t = 0 \\ p(s_{t+1} | s_t, a_t) \delta(b_{t+1}) \pi(a_{t+1} | s_t) & \text{if } b_t = 1 \end{cases}, \quad (20)$$

$$(4) \text{ state-reward function} \quad \bar{\rho}(s, b, a) = r(s, a) b. \quad (21)$$

We can construct $\Omega = (Z, \nu, \kappa, \bar{\mathbf{r}})$, a sub-MRP with interval $n = 2$. Since we always skip the step in which $b = 1$, we only have to define the transition kernel for $b_t = 0$, i.e.

$$\kappa(z_{t+1} | z_t) = \sigma^2(s_{t+1}, b_{t+1}, a_{t+1} | s_t, b_t, a_t) \quad (22)$$

$$= \int_{S \times A} \sigma(s_{t+1}, b_{t+1}, a_{t+1} | s', 1, a') \sigma(s', 1, a' | s_t, 0, a_t) d(s', a') \quad (23)$$

$$= \int_{S \times A} p(s_{t+1} | s', a') \delta(b_{t+1}) \pi(a_{t+1} | s') \delta(s' - s_t) \pi(a' | s_t) d(s', a') \quad (24)$$

$$= \int_A p(s_{t+1} | s_t, a') \delta(b_{t+1}) \pi(a' | s_t) da'. \quad (25)$$

For the state-reward function we have (again only considering $b = 0$)

$$\bar{\mathbf{r}}(s, b, a) = v_\Psi^2(s, b, a) \quad (26)$$

$$= \underbrace{\bar{\rho}(s, 0, a)}_{=0} + \int_{S \times A} \bar{\rho}(s', 1, a') \sigma(s', 1, a' | s, 0, a) d(s', a') \quad (27)$$

$$= \int_{S \times A} r(s', a') \delta(s' - s) \pi(a' | s) d(s', a') \quad (28)$$

$$= \int_A r(s, a') \pi(a' | s) da'. \quad (29)$$

The sub-MRP Ω is already very similar to $TBMRP(E, \pi)$ except for having a larger state-space. To get rid of the b and a state components, we reduce Ω with a state transformation $f(s, b, a) = s$. The reduced MRP has

$$(1) \text{ state space} \quad \{f(z) : z \in Z\} = S, \quad (30)$$

$$(2) \text{ initial distribution} \quad \int_{f^{-1}(s)} \nu(z) dz = \int_{\{s\} \times \{0, 1\} \times A} \mu(s) \delta(b) \delta(a - c) d(s, b, a) = \mu(s), \quad (31)$$

$$(3) \text{ transition kernel} \quad \int_{f^{-1}(s_{t+1})} \kappa(z' | z) dz' \text{ for almost all } z \in f^{-1}(s_t) \quad (32)$$

$$= \int_{\{s_{t+1}\} \times \{0, 1\} \times A} \kappa(z' | z) dz' \text{ for almost all } z \in \{s_t\} \times \{0, 1\} \times A \quad (33)$$

$$= \int_A p(s_{t+1} | s_t, a') \pi(a' | s_t) da' \quad (34)$$

$$(4) \text{ state-reward function} \quad \bar{\mathbf{r}}(z) \text{ for almost all } z \in f^{-1}(s). \quad (35)$$

$$= \int_A r(s, a') \pi(a' | s) da', \quad (36)$$

which is exactly $TBMRP(E, \pi)$. \square

Lemma 1. In a Real-Time Markov Decision Process for the action-value function we have

$$q_{RTMDP(E)}^\pi(s_t, a_t, \mathbf{a}_t) = r(s_t, a_t) + \mathbb{E}_{s_{t+1} \sim p(\cdot | s_t, a_t)} [\mathbb{E}_{\mathbf{a}_{t+1} \sim \pi(\cdot | s_{t+1}, \mathbf{a}_t)} [q_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_t, \mathbf{a}_{t+1})]] \quad (6)$$

Proof. After starting with the definition of the action-value function for an environment $(X, A, \mu, p, r) = RTMDP(E)$ with $E = (S, A, \mu, p, r)$, we separate the transition distribution p into its two constituents p and δ and then, integrate over the Dirac delta.

$$q_{RTMDP(E)}^\pi(s_t, \mathbf{a}_t) = q_{RTMDP(E)}^\pi(s_t, a_t, \mathbf{a}_t) \quad (37)$$

$$= r(s_t, a_t, \mathbf{a}_t) + \mathbb{E}_{s_{t+1}, \mathbf{a}_{t+1} \sim p(\cdot | s_t, a_t, \mathbf{a}_t)} [\underbrace{\mathbb{E}_{\mathbf{a}_{t+1} \sim \pi(\cdot | s_{t+1}, \mathbf{a}_{t+1})} [q_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_{t+1}, \mathbf{a}_{t+1})]}_{(38)}] \quad (38)$$

$$= r(s_t, a_t) + \int_S p(s_{t+1} | s_t, a_t) \int_A \delta(a_{t+1} - \mathbf{a}_t) \dots da_{t+1} ds_{t+1} \quad (39)$$

$$= r(s_t, a_t) + \int_S p(s_{t+1} | s_t, a_t) \mathbb{E}_{\mathbf{a}_{t+1} \sim \pi(\cdot | s_{t+1}, \mathbf{a}_t)} [q_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_t, \mathbf{a}_{t+1})] ds_{t+1} \quad (40)$$

□

Lemma 2. In a Real-Time Markov Decision Process for the state-value function we have

$$v_{RTMDP(E)}^\pi(s_t, a_t) = r(s_t, a_t) + \mathbb{E}_{s_{t+1} \sim p(\cdot | s_t, a_t)} [\mathbb{E}_{\mathbf{a}_t \sim \pi(\cdot | s_t, a_t)} [v_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_t)]] \quad (8)$$

Proof. We follow the same procedure as for Lemma 1.

$$v_{RTMDP(E)}^\pi(\mathbf{x}_t) = v_{RTMDP(E)}^\pi(s_t, a_t) \quad (41)$$

$$= \mathbb{E}_{\mathbf{a}_t \sim \pi(\cdot | s_t, a_t)} [r(s_t, a_t, \mathbf{a}_t) + \mathbb{E}_{s_{t+1}, \mathbf{a}_{t+1} \sim p(\cdot | s_t, a_t, \mathbf{a}_t)} [v_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_{t+1})]] \quad (42)$$

$$= r(s_t, a_t) + \mathbb{E}_{\mathbf{a}_t \sim \pi(\cdot | s_t, a_t)} [\int_S p(s_{t+1} | s_t, a_t) \int_A \delta(a_{t+1} - \mathbf{a}_t) v_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_{t+1}) da_{t+1} ds_{t+1}] \quad (43)$$

$$= r(s_t, a_t) + \int_S p(s_{t+1} | s_t, a_t) \mathbb{E}_{\mathbf{a}_t \sim \pi(\cdot | s_t, a_t)} [v_{RTMDP(E)}^\pi(s_{t+1}, \mathbf{a}_t)] ds_{t+1} \quad (44)$$

□

Proposition 1. The following policy loss based on the state-value function

$$L_{RTMDP(E), \pi}^{RTAC} = \mathbb{E}_{(s_t, a_t) \sim D} \mathbb{E}_{s_{t+1} \sim p(\cdot | s_t, a_t)} D_{KL}(\pi(\cdot | s_t, a_t) || \exp(\frac{1}{\alpha} \gamma v(s_{t+1}, \cdot)) / Z(s_{t+1})) \quad (10)$$

has the same policy gradient as $L_{RTMDP(E), \pi}^{SAC}$, i.e.

$$\nabla_{\pi} L_{RTMDP(E), \pi}^{RTAC} = \nabla_{\pi} L_{RTMDP(E), \pi}^{SAC} \quad (11)$$

Proof. As shown in Haarnoja et al. (2018a), Equation 9 can be reparameterized to obtain the policy gradient, which, applied in a RTMDP, yields

$$\nabla_{\pi} L_{RTMDP(E), \pi}^{SAC} = \mathbb{E}_{\mathbf{x}_t, \epsilon} [\nabla_{\pi} (\log \pi(\mathbf{h}_{\pi}(\mathbf{x}_t, \epsilon), \mathbf{x}_t) - \frac{1}{\alpha} \gamma v(\mathbf{x}_t, \mathbf{h}_{\pi}(\mathbf{x}_t, \epsilon)))] \quad (45)$$

and reparameterizing Equation 10 yields

$$\nabla_{\pi} L_{RTMDP(E), \pi}^{RTAC} = \mathbb{E}_{\mathbf{x}_t, \epsilon} [\nabla_{\pi} (\log \pi(\mathbf{h}_{\pi}(\mathbf{x}_t, \epsilon), \mathbf{x}_t) - \frac{1}{\alpha} \gamma \nabla_{\pi} \mathbb{E}_{s_{t+1} \sim p(\cdot | \mathbf{x}_t)} [v(s_{t+1}, \mathbf{h}_{\pi}(\mathbf{x}_t, \epsilon))])] \quad (46)$$

where \mathbf{h}_{π} is a function mapping from state and noise to an action distributed according to π . This leaves us to show that

$$\nabla_{\mathbf{a}_t} q(\mathbf{x}_t, \mathbf{a}_t) = \underbrace{\nabla_{\mathbf{a}_t} r(\mathbf{x}_t, \mathbf{a}_t)}_{=0} + \nabla_{\mathbf{a}_t} \gamma \mathbb{E}_{s_{t+1} \sim p(\cdot | \mathbf{x}_t, \mathbf{a}_t)} [v(\mathbf{x}_{t+1})] = \gamma \nabla_{\mathbf{a}_t} \mathbb{E}_{s_{t+1} \sim p(\cdot | \mathbf{x}_t)} [v(s_{t+1}, \mathbf{a}_t)] \quad (47)$$

which follows from the definition of the soft action-value function and simplifying quantities defined in the RTMDP. □

D Definitions

Definition 7. A Turn-Based Markov Decision Process $(Z, A, \nu, q, \rho) = \text{TBMDP}(E)$ augments another Markov Decision Process $E = (S, A, \mu, p, r)$, such that

- (1) state space $Z = S \times \{0, 1\}$,
- (2) action space A ,
- (3) initial state distribution $\nu(s_0, b_0) = \mu(s_0) \delta(b_0)$,
- (4) transition distribution $q(s_{t+1}, b_{t+1} | s_t, b_t, a_t) = \begin{cases} \delta(s_{t+1} - s_t) \delta(b_{t+1} - 1) & \text{if } b_t = 0 \\ p(s_{t+1} | s_t, a_t) \delta(b_{t+1}) & \text{if } b_t = 1 \end{cases}$
- (5) reward function $\rho(s, b, a) = r(s, a) b$.

Definition 8. $\Omega = (Z, \nu, \kappa, \bar{r})$ is a sub-MRP of $\Psi = (Z, \nu, \sigma, \bar{\rho})$ if its states are sub-sampled with interval $n \in \mathbb{N}$ and rewards are summed over each interval, i.e. for almost all z

$$\kappa(z' | z) = \kappa^n(z' | z) \quad \text{and} \quad \bar{r}(z) = v_\Psi^n(z). \quad (48)$$

Definition 9. A MRP $\Omega = (S, \mu, \kappa, \bar{r})$ is a reduction of $\Omega = (Z, \nu, \kappa, \bar{r})$ if there is a state transformation $f : Z \rightarrow S$ that neither affects the evolution of states nor the rewards, i.e.

$$(1) \text{ state space} \quad S = \{f(z) : z \in Z\}, \quad (49)$$

$$(2) \text{ initial distribution} \quad \mu(s) = \int_{f^{-1}(s)} \nu(z) dz, \quad (50)$$

$$(3) \text{ transition kernel} \quad \kappa(s_{t+1} | s) = \int_{f^{-1}(s_{t+1})} \kappa(z' | z) dz' \text{ for almost all } z \in f^{-1}(s), \quad (51)$$

$$(4) \text{ state-reward function} \quad r(s) = \bar{r}(z) \text{ for almost all } z \in f^{-1}(s). \quad (52)$$

Definition 10. A MRP Ψ contains another MRP Ω (we write $\Omega \propto \Psi$) if Ψ works at a higher frequency and has a richer state than Ψ but behaves otherwise identically. More precisely,

$$\Omega \propto \Psi \iff \Omega \text{ is a reduction (Def. 9) of a sub-MRP (Def. 8) of } \Psi. \quad (53)$$

Definition 11. The n -step transition function of a MRP $\Omega = (S, \mu, \kappa, \bar{r})$ is

$$\kappa^n(s_{t+n} | s_t) = \int_S \kappa(s_{t+n} | s_{t+n-1}) \kappa^{n-1}(s_{t+n-1} | s_t) ds_{t+n-1}. \quad | \text{ with } \kappa^1 = \kappa \quad (54)$$

Definition 12. The n -step value function v_Ω^n of a MRP $\Omega = (S, \mu, \kappa, \bar{r})$ is

$$v_\Omega^n(s_t) = \bar{r}(s_t) + \int_S \kappa(s_{t+1} | s_t) v_\Omega^{n-1}(s_{t+1}) ds_{t+1}. \quad | \text{ with } v_\Omega^1 = \bar{r} \quad (55)$$