

# Dynamics of innovation along foraged trajectories

S Ganga Prasath

## 1 2D grid world

Let us say an agent has an intrinsic behavior to follow a herd or a trail laid by other ants or a flock of birds. All the agents simply following this herd following behavior will not result in new solutions. Agents often have to innovate and find new solutions either because the trails are not available anymore or because you know a better route (influence of history) or perturbations (intrinsic or environmental) can throw you off trails and you have to find new solutions. We are interested in understanding the role of intrinsic behavior on the dynamics of innovation.

We will start by implementing an agent trying to optimally traverse from point A to point B in a 2D grid world. The steps involved in traditional reinforcement learning based on SARSA algorithm involves essentially 2-steps: (i) action-value function update, (ii) policy update. The action-value update equation in on-policy SARSA is given by

$$Q_{\pi}(s_t, a_t) \leftarrow Q_{\pi}(s_t, a_t) + \alpha \{r_{t+1} + \gamma Q_{\pi}(s_{t+1}, a_{t+1}) - Q_{\pi}(s_t, a_t)\} \quad (1)$$

while the policy update is given by

$$\pi(s) = \arg \max_a Q(s, a).$$

This is shown in algorithmic form in Alg. 1.

### Algorithm 1: SARSA algorithm for action-value update

```

Initialize: State,  $s_0$ ; action,  $a_0$ ; action-value function,  $Q(s_j, a_j)$ 
foreach epoch do
    Update action-value function using SARSA rule:

         $Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) + \alpha\{r_{t+1} + \gamma Q_\pi(s_{t+1}, a_{t+1}) - Q_\pi(s_t, a_t)\}$ 

    Choose action:

         $\pi(s_t) = \arg \max_a Q(s_t, a_t)$ 

    Calculate reward,  $r_t$ 
    Check if target  $s^*$  is reached, continue if not
end

```

**Algorithm 1:**

## 2 Trail following behavior

Consider a trail of pheromone  $\mathbf{x}^*(s)$  represented by arc-length parameterization  $s$ . The concentration field in 2D is then given by  $c(\mathbf{x}) = c_o \delta(\mathbf{x} - \mathbf{x}^*(s))$ . This field of course is assumed to be steady but can be made time-dependent by simply adding a decay time-scale  $\tau$  to get:  $c(\mathbf{x}, t) = c_o \delta(\mathbf{x} - \mathbf{x}^*(s)) e^{-t/\tau}$ . We can immediately evaluate some of the properties of the curve  $\mathbf{x}^*(s)$  which will come in handy soon:  $\hat{\mathbf{t}}(s) = d\mathbf{x}^*(s)/ds = \{\cos \psi(s), \sin \psi(s)\}$  and  $\hat{\mathbf{n}}(s) = \{\sin \psi(s), \cos \psi(s)\}$ . As we can see the entire curve  $\mathbf{x}^*(s)$  can be represented only using  $\psi(s)$  up to global translations and rotations, which is a well known property of curves in 2D.

The state of the agent/ant in our problem is represented by its coordinates  $\mathbf{r}(t) = \{r_x(t), r_y(t)\}$  and it can make measurements about how far from the pheromone trail it is  $d\mathbf{r}(t)$  as well as the orientation of the trail,  $\hat{\mathbf{t}}(s, t)$ . The action that it takes from these measurements is to align its orientation,  $\hat{\mathbf{p}}(t)$  along a direction that will take it towards the trail. We show in Fig. 2 a schematic of the problem set up. The strategy used by the agent for tracking the pheromone trail will be to move along the direction given by  $(d\mathbf{r}/|d\mathbf{r}| + \hat{\mathbf{t}})$  by a fixed length  $h$ . In order to identify the location along  $\mathbf{x}^*(s)$  where  $d\mathbf{r}$  intersects, we use the condition  $d\mathbf{r}(t) \perp \hat{\mathbf{t}}(s, t)$ .

### 2.1 Semi-circular trail

For the problem at hand, we will consider a semi-circular trail whose coordinates can be written in arc-length form as  $\mathbf{x}^*(s) = \{a \cos(s/a), a \sin(s/a)\}$  where  $a$  is the radius of the circle. The trail starts at  $s = 0$  given by coordinates  $\mathbf{x}^*(0) = \{a, 0\}$  and ends at  $s = \pi a$

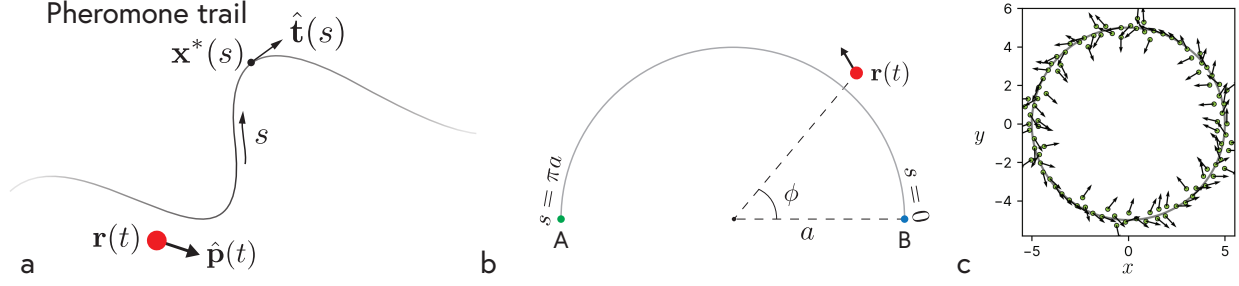


Figure 1: Schematic of setup

at  $\mathbf{x}^*(\pi a) = \{-a, 0\}$ . From this it is easy to find that  $\hat{\mathbf{t}}(s) = \{-\sin(s/a), \cos(s/a)\} = \{\cos(\pi/2 + s/a), \sin(\pi/2 + s/a)\}$ , and  $\hat{\mathbf{n}}(s) = \{-\cos(s/a), \sin(s/a)\}$ . We see that we can represent  $\hat{\mathbf{t}}(s)$  through  $\psi(s) = (\pi/2 + s/a)$ . We can now calculate the location along the arc-length where  $\mathbf{r}(t)$  is closest by using the constraint condition  $d\mathbf{r}(t) \perp \hat{\mathbf{t}}(s, t)$  or equivalently  $d\mathbf{r}(t) \cdot \hat{\mathbf{t}}(s, t) = 0$ . We can now write  $d\mathbf{r}(t) \sim \{a \cos(s/a) - r_x, a \sin(s/a) - r_y\}$  (up to normalization) and get the location  $s_*(t) = a \tan^{-1}(r_y/r_x)$  by using the formula for  $\hat{\mathbf{t}}(s)$ . From this it is trivial to see that the angle along  $d\mathbf{r}$  is  $\phi(t) = \tan^{-1}(r_y/r_x)$  (as is evident in Fig. 2(b)). For a given location of the agent,  $\mathbf{r}(t)$  the orientation it needs to take in the next step can be easily calculated to be  $\vartheta(t) = (\pi/4 - \phi(t))$ .

We can now set up the entire dynamics of the agent's trail following behavior on this semi-circle. We can state this in the notation of reinforcement learning as it will become helpful later. The state of the agent is  $S^t = \{r_x^t, r_y^t\}$ , the measurements it makes are  $M^t = \phi^t$  and using this information the action the agent takes is  $A^t = \{\vartheta^t\}$  which is its orientation. The dynamics of the agent can now be written as follows:

$$\text{Measurement update: } \phi^{t+1} = \tan^{-1} \left( \frac{r_y^t}{r_x^t} \right) + \zeta^t, \quad (2)$$

$$\text{Action update: } \hat{\mathbf{p}}^{t+1} = \hat{\mathbf{t}}^{t+1} + d\mathbf{r}^{t+1}, \quad (3)$$

$$\text{State update: } \mathbf{r}^{t+1} = \mathbf{r}^t + l\hat{\mathbf{p}}^{t+1}, \quad (4)$$

where  $\mathbf{r}^t = \{r_x^t, r_y^t\}$ ,  $d\mathbf{r}^t = (\mathbf{x}^*(s^t) - \mathbf{r}^t)/\|\mathbf{x}^*(s^t) - \mathbf{r}^t\|$  and  $\hat{\mathbf{p}}^t = \{\cos \vartheta^t, \sin \vartheta^t\}$ . We have added sensory noise  $\zeta^t$  (which is sampled from a uniform distribution) to the measurement to reflect the error that accompanies measurements usually. The solution dynamics is shown in Fig. 2(c). We call this line-following behavior as a policy  $\pi_{\text{ite}}(\mathbf{r}, \phi)$  : which denotes the agent's innate behavior to follow pheromone trails.

### 3 Ornstein-Uhlenbeck process

Let us consider an agent whose orientation  $\hat{\mathbf{p}} = (\cos \vartheta, \sin \vartheta)$  wants to relax to a neutral orientation with stochasticity associated with estimation error. The dynamics of the orientation

can be written as a stochastic differential equation given by

$$\dot{\vartheta}(t) = -\alpha\vartheta(t) + \sqrt{2D}\eta(t), \quad (5)$$

$$\langle \eta(t)\eta(0) \rangle = \delta(t). \quad (6)$$

This is the evolution of the famous Ornstein-Uhlenbeck process which has been extremely well studied and the evolution of the probability density can be solved exactly. The Fokker-Planck equation for its dynamics which provides the probability of a given state  $\vartheta(t)$  i.e.  $p(\vartheta, t)$ , is given by

$$\frac{\partial \mathcal{P}}{\partial t} = \alpha \frac{\partial(\vartheta \mathcal{P})}{\partial \vartheta} + D \frac{\partial^2 \mathcal{P}}{\partial \vartheta^2}$$

The dynamics of the angle relaxes to equilibrium angle  $\vartheta = 0$  in the absence of fluctuations. The evolution of the mean is  $\dot{h}(t) = -\alpha h(t)$  where  $h(t) = \mathbb{E}[\vartheta(t)]$  while the co-variance dynamics is given by  $\dot{\zeta}(t) = -2\alpha\zeta(t) + 2D$  with  $\zeta = \mathbb{E}[(\vartheta(t) - h(t))^2] = \mathbb{E}[\vartheta^2(t)] - h(t)^2$ .

The position of the agent evolves as  $\dot{\mathbf{r}}(t) = v_o \hat{\mathbf{p}}(t)$ . Under the assumption that the angle is small, we can write  $\mathbf{r}(t) = v_o(1, \vartheta)$ . Thus we have  $\dot{y}(t) = v_o\vartheta$  which can be used to rewrite the orientation equation as

$$\ddot{y}(t) = -\alpha\dot{y}(t) + v_o\sqrt{2D}\eta(t). \quad (7)$$

This second order equation needs to boundary condition. We require  $y(0) = 0$  and  $\dot{y}(0) = 0$  which comes from setting the initial orientation to be 0. We are interested in the evolution of the mean and the variance of this SDE. Let  $m(t) = \mathbb{E}[y(t)]$  where the averaging is the ensemble average (equivalently the average of the several trajectories of the same dynamics). The mean  $m(t)$  evolves as

$$\ddot{m}(t) = -\alpha\dot{m}(t). \quad (8)$$

By applying the above boundary condition, we can immediately see that  $m(t) = 0$ . Nevertheless, we are interested in the evolution of the standard deviation,  $\beta(t) = \mathbb{E}[(y(t) - m(t))^2]$ . In order to derive an expression for that, we differentiate the definition once to get

$$\begin{aligned} \dot{\beta}(t) &= 2\mathbb{E}[\{y(t) - m(t)\}\{\dot{y}(t) - \dot{m}(t)\}], \\ &= 2\mathbb{E}[y\dot{y}] - 2m\dot{m}, \\ &= 2v_o\{\mathbb{E}[y\vartheta] - m\dot{m}\}. \end{aligned}$$

We can estimate  $\mathbb{E}[y\vartheta]$  as follows

$$\begin{aligned} \mathbb{E}[y\vartheta] &= v_o\mathbb{E}\left[\int_0^t \vartheta(s)\vartheta(t) \, ds\right], \\ &= v_o\int_0^t \mathbb{E}[\vartheta(t)\vartheta(s)] \, ds. \end{aligned}$$

We know that for an OU process,

$$\mathbb{E}[\vartheta(t)\vartheta(s)] = \frac{D}{\alpha} \left[ e^{-\alpha(t-s)} - e^{-\alpha(t+s)} \right].$$

From this we see that

$$\begin{aligned} \mathbb{E}[y\vartheta] &= \frac{v_o D}{\alpha} \int_0^t \left[ e^{-\alpha(t-s)} - e^{-\alpha(t+s)} \right] ds, \\ &= \frac{v_o D}{\alpha^2} \left[ 1 - 2e^{-\alpha t} + e^{-2\alpha t} \right]. \end{aligned}$$

We thus have the evolution of the co-variance of  $y(t)$  as

$$\dot{\beta}(t) = \frac{2v_o^2 D}{\alpha^2} \left[ 1 - 2e^{-\alpha t} + e^{-2\alpha t} \right]. \quad (9)$$

We finally have the variance as a function of time as

$$\beta(t) = \frac{v_o^2 D}{\alpha^3} \left[ 2\alpha t - e^{-2\alpha t} + 4e^{-\alpha t} - 3 \right]. \quad (10)$$

For  $t \gg \alpha^{-1}$  we can immediately see that  $\beta(t) \approx 2v_o^2 D t / \alpha^2$ . We are interested in the variance at time,  $t^* = c/v_o$  where  $c$  is the distance of the agent from the goal A and  $v_o$  is the intrinsic speed of the agent. At this time the location of the agent along x-axis is  $x(t^*) = c$  while the variance in  $y(t)$  is  $\beta(t^*) \approx 2v_o D c / \alpha^2$ .

## 4 Evaluating value function of policies

We are at a stage where we can solve the speed/distance-accuracy trade-off for the an agent trying to make optimal decision on when to initiate innovation. In order to do that, let us define clearly what is the state and action of the agent.

### 4.1 $v(\mathbf{r}, \vartheta)$ for $\pi_{\text{int}}(\mathbf{r}, \vartheta)$

The agent's state  $s^t$  is given by its position,  $s^t : \mathbf{r}^t = \{x^t, y^t\}$  and the action it can take at these location is choose a orientation to move,  $a^t : \vartheta^t$ . Once the agent choose this action, it is taken to the next state  $s^t$  through the dynamics:  $\mathbf{r}^{t+1} = \mathbf{r}^t + l\hat{\mathbf{p}}_{\vartheta}^t$ . Before we compute the state-value function for the intrinsic policy,  $\pi_{\text{int}}(s)$ , the sequence of process we described above can be written as

$$\text{Measurement: } \phi^t = \tan^{-1} \left( \frac{r_y^t}{r_x^t} \right), \quad (11)$$

$$\text{Action, } a^t: \vartheta^t = \frac{\pi}{2} + \phi^t + \zeta^t, \quad (12)$$

$$\text{State update, } s^{t+1}: \mathbf{r}^{t+1} = \mathbf{r}^t + l\hat{\mathbf{p}}_{\vartheta}^t. \quad (13)$$

It is however difficult to evaluate the value function for this specific functional form of the policy. We assume that the distribution of the orientation of the agent is distributed as through it is a worm-like chain polymer whose orientation distribution is given by

$$\mathcal{P}[\vartheta(s)] = \frac{1}{\mathcal{Z}} \exp(-\beta E_b)$$

where the bending energy  $E_b[\vartheta(s)] = B/2 \int_0^L (\vartheta'(s) - \kappa_o)^2 ds$ , fugacity  $\mathcal{Z} = \int \mathcal{D}[\vartheta(s)] \exp\{-\beta E_b[\vartheta(s)]\}$  and the bending stiffness  $B$  captures the constraint for the agent to move forward,  $\beta$  is the effective temperature capturing the fluctuations in the trajectory of agent. We also have boundary conditions  $\vartheta(0) = \pi/2$  and  $\vartheta'(0) = \kappa_o = a^{-1}$ .

We are interested in evaluating the value function of this policy,  $v(s^t)$  whose expression is given by

$$v_{\text{int}}(\mathbf{r}, \vartheta) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \middle| S_t = S \right].$$

We can write this in path-integral form as

$$v_{\text{int}}(\mathbf{r}, \vartheta) = \sum_{k=t}^T \int \mathcal{D}[\vartheta(s)] \mathcal{P}[\vartheta(s)] \gamma^{k-t} R_{k+1}, \quad (14)$$

$$\approx \int_s^L \int \mathcal{D}[\vartheta(s)] \mathcal{P}[\vartheta(s)] R(s) ds \quad (15)$$

The reward itself is a function of the agent's azimuthal direction,  $\phi(s) = \tan^{-1}(y(s)/x(s))$ . Further we know the geometric relation  $x'(s) = \cos(\vartheta(s))$ ,  $y'(s) = \sin(\vartheta(s))$ . The functional form of the reward is  $R(\phi) = \exp(-\phi/\phi^*)$  and

$$\phi(\mathbf{r}, \vartheta) = \tan^{-1} \left( \frac{\int_0^s \sin(\vartheta(\tau)) d\tau}{\int_0^s \cos(\vartheta(\tau)) d\tau} \right)$$

We can write the ultimate expression for the value function we would like to evaluate as,

$$v_{\text{int}}(\mathbf{r}, \vartheta) = \int_s^L \int \mathcal{D}[\vartheta(s)] \mathcal{P}[\vartheta(s)] \exp(-\phi(s)/\phi^*) ds'. \quad (16)$$

We know that  $\vartheta(s)$  has to satisfy the boundary conditions  $\vartheta(0) = \pi/2$  and  $\vartheta'(0) = \kappa_o$ . Before we solve the stochastic version of the problem, let us solve it for the deterministic case when the agent follows the semi-circular path exactly. In this scenario,  $\vartheta_b(s) = \pi/2 + \kappa_o s$  and from this it immediately follows that  $\phi_b(s) = \kappa_o s$ . For this particular functional form of  $\phi(s)$  we can write the reward as  $R(s) = \exp(-\kappa_o s/\phi^*)$ . The value function from this can be written as

$$v_b(s) = \int_s^L \exp(-\kappa_o \tau/\phi^*) d\tau, \quad (17)$$

$$= -\frac{\phi^*}{\kappa_o} [e^{-\kappa_o L/\phi^*} - e^{-\kappa_o s/\phi^*}], \quad (18)$$

$$= \frac{\phi^*}{\kappa_o} e^{-\kappa_o L/\phi^*} [e^{\kappa_o (L-s)/\phi^*} - 1]. \quad (19)$$

We now introduce perturbations in the trajectory of the agent from the trajectory  $\vartheta_b(s)$  by introducing perturbations  $\vartheta(s) = \vartheta_b(s) + \delta\vartheta(s)$ . We know that  $\phi$  and  $\vartheta$  are related by the geometric constraint:  $\phi(s) = \vartheta(s) - \pi/2$  which immediately gives  $\phi(s) = \kappa_o s + \delta\vartheta(s)$ . We can then write the value function as

$$v_{\text{int}}(s) = \int_s^L \int \frac{\mathcal{D}[\vartheta(s)]}{\mathcal{Z}} e^{-(\beta B/2) \int_0^L (\vartheta'(s) - \kappa_o)^2 ds} \exp(-\phi(\tau)/\phi^*) d\tau, \quad (20)$$

$$= \int_s^L \int \frac{\mathcal{D}[\vartheta(s)]}{\mathcal{Z}} e^{-(\beta B/2) \int_0^L [\delta\vartheta'(s)]^2 ds} e^{(-\kappa_o \tau - \delta\vartheta(\tau))/\phi^*} d\tau. \quad (21)$$

Let us start by calculating  $\mathcal{Z}$ ,

$$\mathcal{Z} = \int \mathcal{D}[\vartheta(s)] \exp\{-\beta E_b[\vartheta(s)]\}, \quad (22)$$

$$= \int \mathcal{D}[\vartheta(s)] e^{-(\beta B/2) \int_0^L [\delta\vartheta'(s)]^2 ds}, \quad (23)$$

$$\delta\vartheta(s) = \sum_{k=-\infty}^{\infty} \hat{a}_k e^{i2\pi ks/L}, \quad (24)$$

$$\int_0^L [\delta\vartheta'(s)]^2 ds = \int_0^L \left[ \sum_{k=-\infty}^{\infty} \frac{2\pi k}{L} \hat{a}_k i e^{i2\pi ks/L} \right] \left[ \sum_{q=-\infty}^{\infty} \frac{2\pi q}{L} \hat{a}_q i e^{i2\pi qs/L} \right] ds, \quad (25)$$

$$= \frac{4\pi^2}{L} \sum_{k=-\infty}^{\infty} k^2 \hat{a}_k^2, \text{ where we have used } \hat{a}_{-k} = \hat{a}_k \quad (26)$$

$$\mathcal{Z} = \prod_{k=-\infty}^{\infty} \int_{-\infty}^{\infty} d\hat{a}_k e^{-\frac{\beta B k^2}{L} 4\pi^2 \hat{a}_k^2}, \quad (27)$$

$$= \prod_{k=-\infty}^{\infty} \sqrt{\frac{\pi L}{\beta B}} \frac{1}{2\pi k} \quad (28)$$

We can go ahead and write down the value function as

$$v_{\text{int}}(s) = \frac{1}{\mathcal{Z}} \prod_{k=-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{\beta B k^2}{L} 4\pi^2 \hat{a}_k^2} \left[ \int_s^L e^{(-\kappa_o \tau - \hat{a}_k e^{i2\pi k\tau/L})/\phi^*} d\tau \right] d\hat{a}_k, \quad (29)$$

$$= \frac{1}{\mathcal{Z}} \prod_{k=-\infty}^{\infty} \int_s^L \left[ \int_{-\infty}^{\infty} e^{-\frac{\beta B k^2}{L} 4\pi^2 \hat{a}_k^2 - \hat{a}_k e^{i2\pi k\tau/L}/\phi^*} d\hat{a}_k \right] e^{-\kappa_o \tau/\phi^*} d\tau. \quad (30)$$

We can evaluate the integral inside the square brackets as

$$\int_{-\infty}^{\infty} e^{-\alpha \hat{a}^2 - \gamma \hat{a}} d\hat{a} = e^{-\frac{\gamma^2}{4\alpha}} \int_{-\infty}^{\infty} e^{-\alpha (\hat{a} + \frac{\gamma}{2\alpha})^2} d\hat{a} = \sqrt{\frac{\pi}{\alpha}} e^{-\frac{\gamma^2}{4\alpha}}.$$

We then get

$$v_{\text{int}}(s) = \sqrt{\frac{\pi L}{\beta B}} \frac{1}{\mathcal{Z}} \prod_{k=-\infty}^{\infty} \frac{1}{2\pi k} \int_s^L e^{-\gamma(k,\tau)^2/4\alpha} e^{-\kappa_o \tau/\phi^*} d\tau. \quad (31)$$

where  $\gamma(k, \tau) = e^{i2\pi k \tau/L}/\phi^*$ ,  $\alpha = \beta B k^2 4\pi^2/L$ . In the limit of small temperature, or large bending stiffness,  $\alpha \gg 1$  and we can Taylor expand the integral to get

$$v_{\text{int}}(s) = \sqrt{\frac{\pi L}{\beta B}} \frac{1}{\mathcal{Z}} \prod_{k=-\infty}^{\infty} \frac{1}{2\pi k} \int_s^L \left[ 1 - \frac{\gamma(k, \tau)^2}{4\alpha} \right] e^{-\kappa_o \tau/\phi^*} d\tau. \quad (32)$$

From this we see that the leading order contribution is still

$$v_{\text{int}}(s) = \frac{\phi^*}{\kappa_o} e^{-\kappa_o L/\phi^*} [e^{\kappa_o(L-s)/\phi^*} - 1]. \quad (33)$$

## 4.2 $v(\mathbf{r}, \vartheta)$ for $\pi_{\text{OU}}(\mathbf{r}, \vartheta)$

We have seen already the dynamics of the OU process and we can write the policy  $\pi_{\text{OU}}(\mathbf{r}, \vartheta)$  using the evolution of the orientation and position as

$$\vartheta^{t+1} = \vartheta^t(1 - \mu\Delta t) + \sqrt{2D\Delta t}\zeta^t, \quad (34)$$

$$\langle \zeta^t \zeta^0 \rangle = \delta(t), \quad (35)$$

$$\mathbf{r}^{t+1} = \mathbf{r}^t + l\hat{\mathbf{p}}_{\vartheta}^t, \text{ where } \hat{\mathbf{p}}_{\vartheta}^t = \{\cos \vartheta^t, \sin \vartheta^t\}. \quad (36)$$

We have calculated the evolution of variance of the agent's trajectory from a straight line as a function of distance/time from the bifurcation point evolves as,

$$\psi(t) = \frac{v_o^2 D}{\mu^3} \left[ 2\mu t - e^{-2\mu t} + 4e^{-\mu t} - 3 \right]. \quad (37)$$

where  $\mu^{-1}$  is the time-scale of relaxation of the agent. The  $n$ -step variance of the agent is  $\psi(t^*)$ , where  $t^* = n\Delta t$ . Assuming that the agent reaching a distance  $\sigma$  from the end-point A can reach the target, the probability that an agent will reach the target is given by  $\mathcal{P}_{\text{OU}}(\|\mathbf{r}^{n\Delta t} - \mathbf{r}^*\| \leq \sigma) = \sigma/\psi(t^*)$ . We can now calculate the value function,  $v_{\text{OU}}(\mathbf{r}, \vartheta)$  using the formula

$$v_{\text{OU}}(\mathbf{r}, \vartheta) = \int_0^{t^*} \int \mathcal{D}[\vartheta(t)] \mathcal{P}_{\text{OU}}[\vartheta(t)] R(t) dt.$$

Let us assume that the reward function for the OU process is  $R(t) = \nu[1 - \Theta(\|\mathbf{r}(t^*) - \mathbf{r}^*\| - \sigma)]$ . From this we see that the value of  $v_{\text{OU}}(\mathbf{r}^c, \vartheta^c) = \sigma\nu/\psi(t^*)$ .



## 5 Value function of combined policy

The total value function for the agent's policy is  $v(\mathbf{r}, \vartheta) = v_{\text{int}}(\mathbf{r}, \vartheta) + v_{\text{OU}}(\mathbf{r}^c, \vartheta^c)$ . The agent is trying to maximize this value function from a given state,  $\mathbf{r}^t, \vartheta^t$ . The final expression for the value function is

$$v(\mathbf{r}, \vartheta) = \frac{\sigma\nu}{\psi(t^* + t_f)} + \frac{\phi^*}{\kappa_o} e^{-\kappa_o s^* / \phi^*} [e^{\kappa_o(s^* - s) / \phi^*} - 1],$$

where  $\psi(t) = \frac{v_o^2 D}{\mu^3} \left[ 2\mu t - e^{-2\mu t} + 4e^{-\mu t} - 3 \right]$ ,  $t_f$  is the time-taken to reach the region around the target location A. In this equation, the diffusion coefficient representing the fluctuations in the dynamics of the agent capturing the error in memory of target location A. We assume that the fluctuation  $D$  is larger farther from the target location.