

---

# EXPLORATION OF THE APPLICABILITY OF PROBABILISTIC INFERENCE FOR LEARNING CONTROL IN UNDERACTUATED AUTONOMOUS UNDERWATER VEHICLES

---

A PREPRINT

**Wilmer Ariza Ramirez**  
Australian Maritime College  
University of Tasmania  
Newnham, Australia  
Wilmer.ArizaRamirez@utas.edu.au

**Zhi Q. Leong**  
Australian Maritime College  
University of Tasmania  
Newnham, Australia  
Zhi.Leong@utas.edu.au

**H. Nguyen**  
Australian Maritime College  
University of Tasmania  
Newnham, Australia  
h.d.nguyen@utas.edu.au

**S.G. Jayasinghe**  
Australian Maritime College  
University of Tasmania  
Newnham, Australia  
shantha.jayasinghe@utas.edu.au

December 30, 2019

## ABSTRACT

Underwater vehicles are employed in the exploration of dynamic environments where tuning of a specific controller for each task would be time-consuming and unreliable as the controller depends on calculated mathematical coefficients in idealised conditions. For such a case, learning task from experience can be a useful alternative. This paper explores the capability of probabilistic inference learning to control autonomous underwater vehicles that can be used for different tasks without re-programming the controller. Probabilistic inference learning uses a Gaussian process model of the real vehicle to learn the correct policy with a small number of real field experiments. The use of probabilistic reinforced learning looks for a simple implementation of controllers without the burden of coefficients calculation, controller tuning or system identification. A series of computational simulations were employed to test the applicability of model-based reinforced learning in underwater vehicles. Three simulation scenarios were evaluated: waypoint tracking, depth control and 3D path tracking control. The 3D path tracking is done by coupling together a line-of-sight law with probabilistic inference for learning control. As a comparison study LOS-PILCO algorithm can perform better than a robust LOS-PID. The results shows that probabilistic model based reinforced learning is a possible solution to motion control of underactuated AUVs as can generate capable policies with minimum quantity of episodes.

**Keywords** PILCO · LOS · Underwater Vehicle · Path tracking · Reinforced learning

## Acknowledgements

The authors thank Defence Science and Technology Group for the loan of the vehicle MULLAYA to the Australian Maritime College, and constant support on the platform development.

## 1 Introduction

Autonomous underwater vehicles play an important role in the exploration of the seas. This exploration is primarily driven by commercial, military and scientific needs. In this context, the proper and correct navigation of the vehicle is a key requirement. Motion controllers that are used for navigating AUVs can be classified in four basic strategies: point stabilization [1], trajectory tracking [2], path following [3] and path tracking [4]. Point stabilization controllers stabilize a vehicle to the desired goal posture from an initial configuration [5]. Trajectory tracking controllers use a virtual vehicle to generate a reference trajectory that has an associated time required to be employed by the real vehicle [6]. In the case of path following the vehicle is forced to pursue the desired path without temporal specifications. In the case of path following controller, they usually employ the Frenet-Serret line-of-sight(LOS) coupled with another controller to minimize the error between the obtained geometric references and the vehicle variables. The final strategy is path tracking, which combines trajectory tracking and path following by the introduction of a virtual time parameter to force the vehicle to complete the path within a specific time.

Difficulty of controlling underwater vehicles arises due to non-linear and time-varying dynamics of underwater vehicles, uncertainties in its hydrodynamic coefficients and disturbance in the environment.(e.g. ocean currents). Furthermore, all complexities are exacerbated for the controller in underactuated vehicles [5]; underactuated vehicles have more degrees of freedom to be controlled than surfaces of control. Nevertheless, this configuration is more prevalent as it is the most energy efficient design for travelling at high speeds [7].

Waypoint tracking is the most common methodology to control a vehicle, e.g. commercial vehicles such as Gavia [8]and REMUS [9] use this methodology. Waypoint tracking is directing the vehicle to approximate to a series of specific target points. The vehicle calculates the required direction to which the vehicle should be directed and upon arriving at the proximities of a point is given a new target. Like many industries, PID is the most common methodology to control the vehicle orientation and speed. However, there is research to employ more robust options than PID. [10] employs a NARMAX model from the vehicle and a constrained self-tuning controller to direct the vehicle to the respective target. Other methodologies employ a combination of LOS for waypoint [11] and a standard controller or backstepping techniques to minimize the error between the vehicle position and the desired position [12].

In the case of LOS, commercially LOS-PID and LOS-Fuzzy controllers are employed as their implementation are simpler and more accepted in the industry [13, 14, 15]. Other research as [16] used a LOS guidance law with two integrators and three feedback controllers to compensate for external unknown perturbation such as ocean current which is one of the weaknesses for methodologies as PID/Fuzzy controllers. [17] have used an alternative methodology of a grey prediction to obtain the next AUV position in advance and then use LOS to calculate the desired angles, such that if there is environmental interference, the vehicle will not be affected.

In the search of more robust controllers, nonlinear control techniques had been explored. [18] have designed a horizontal path following controller based on Lyapunov stability theorem and backstepping method. In [19], a method consisted of Lyapunov stability theorem and feedback gain backstepping reduce the complexity of the controller and improve adjustability of the parameters. Another methodology proposes a global path following for AUV based on the same coordinates to achieve global asymptotic stability of the following error [20]. Following the research of backstepping, [21] have adopted fuzzy backstepping sliding mode control to overcome non-linearities, uncertainties and external disturbances.

However, the aforementioned research in controls are focused to provide a more robust path following performance. The controller still requires the calibration of parameters or specific design of observers to identify the unknown parameters of the dynamic model. A methodology to overcome this is the use of machine learning algorithms. The most prominent algorithm in machine learning is neural networks. In particular, the research to control underwater vehicles had focused on the use of machine learning algorithms to recognize uncertainties. [22, 23] have designed a combined version of control law for the convergence of the kinematic model and an adaptive backstepping sliding control based in radial basis function (RBF) neural network to identify the unknown parameters of the dynamic model. In a similar way, [24] have reduce the backstepping complexity by the inclusion of a second-order filter to obtain the derivatives of the virtual controller and filter high-frequency measurement noise, and coupled the filter with an RBF neural network that compensates for vehicle uncertainties.

Although the practicality of machine learning has been largely to identify uncertainties in AUV control, some machine learning algorithms are capable of doing more such as controlling the vehicle directly. Recently, there has been increasing research efforts on the use of reinforcement learning to generate policies to control underwater vehicles and robots in general. An example of machine learning control can be seen in [25],where reinforced learning (RL) based on the Markov decision process (MDP) was employed to produce a policy capable of controlling a vehicle around an obstacle with a minimum cost. In the case of path-following, reinforced learning had been applied to path following of ships. In [26], an actor-critic multilayer perception reinforced learning is used to reduce the tracking

error to zero. Deep reinforced learning has also been proposed as a possible solution for the tracking problem. [27] employed two neural networks. The primary neural network selects the action and the secondary evaluates whether the produced action is valid; with further modification through a deep deterministic policy gradient. Another application used continuous actor-critic learning automaton algorithm to teach an AUV to follow a pipeline [28], considering the improved performance in search of the policy of this algorithm the number of episodes over the platform can be over the hundreds.

RL can be divided into two methodologies: model-based methods and model-free methods, such as Q-learning [29] or TD-learning [?]. The application of path-following control based in traditional RL such as Q-learning [30] is highly complex and difficult as a high quantity of experiments is required to acquire data and test each policy iteration. The additional difficulties in underwater vehicles are vehicle safety, maximum time underwater and computational power. RL for a system with low-dimensional state spaces and fairly favourable dynamics can require thousands of trials to arrive at the appropriate policies[31, 27].

Model-based RL methods are more efficient than model-free methods in searching for a useful policy, as the policy is searched over a model and not the real platform. However, their accuracy can suffer severely from model errors. A solution to address the model errors is the use of probabilistic models to express its uncertainty. An application of model-free methodologies that use a probabilistic methodology was propose by [32]. Their methodology use a on-line selective reinforcement learning approach combined with Gaussian Process (GP) regression for learning reference tracking control policies given no prior knowledge of the dynamical system. [31] have proposed Probabilistic Inference for Learning Control (PILCO), which is a model-based policy search method. The probabilistic model uses non-parametric Gaussian processes (GPs) to characterise the model uncertainty and the policy improvement is based on analytic policy gradients which employs deterministic approximate interference techniques. Due to probabilistic modelling and inference approach, PILCO can achieve higher learning efficiency than other methods in continuous state-action domains and, hence, is directly applicable to complex mechanical systems, such as robots.

In this paper, the authors explore the applicability of PILCO to control underactuated AUVs by a series of simulations with different objectives and target values. The main goals of our implementation of reinforced learning with PILCO are:

- Minimum quantities of episodes over the platform;
- Small test time over the platform;
- Minimum quantity of variables to be predicted by the GP; and
- Vehicle safety.

## 2 Underwater Vehicle Mathematical Model

In [33] it was shown that the non-linear dynamic equations of motion of an underwater vehicle can be expressed in vector notation defined by a state vector composed by the vector  $v$  of velocities on the body frame of the form  $[u, v, w, p, q, r]^T$  and the vector  $\eta$  of position in the Earth fixed frame (Fig. 1) of the form  $[\xi, \eta, \zeta, \phi, \theta, \psi]^T$  such that

$$\mathbf{M}\dot{\mathbf{v}} + \mathbf{C}(\mathbf{v})\mathbf{v} + \mathbf{D}(\mathbf{v})\mathbf{v} + \mathbf{g}(\eta) = \tau \quad (1)$$

with the kinematic equation

$$\dot{\eta} = \mathbf{J}(\eta)\mathbf{v} \quad (2)$$

where

$\eta$  position and orientation of the vehicle in Earth-fixed frame,

$\mathbf{v}$  linear and angular vehicle velocity in body fixed frame,

$\dot{\mathbf{v}}$  linear and angular vehicle acceleration in body fixed frame,

$\mathbf{M}$  matrix of inertial terms,

$\mathbf{C}(\mathbf{v})$  matrix of Coriolis and centripetal terms,

$\mathbf{D}(\mathbf{v})$  matrix consisting of damping or drag terms,

$\mathbf{g}(\eta)$  vector of restoring forces and moments due to gravity and buoyancy,

$\tau$  vector of control and external forces, and

$\mathbf{J}(\eta)$  rotation matrix that converts velocity in a body fixed frame  $v$  to an Earth fixed frame velocity  $\dot{\eta}$ .

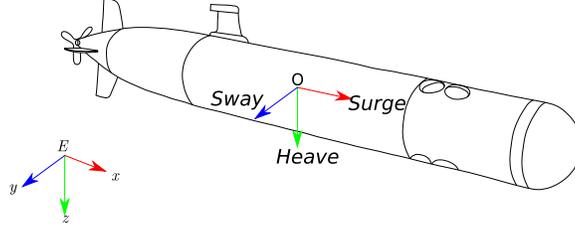


Figure 1: AUV different reference frames, vehicle frame is equal to the centre of buoyancy.

Equation (1) can be expanded into a more general equation of motion as has been shown in [34, 35]. The result of the expansion will be a system of six equation with 73 hydrodynamic coefficients. However, for a complete model the control surfaces must be modelled. In a general case, the resulting forces and moments of a control surface (thrusters and fins) can be expressed as [33]

$$\begin{aligned} F_{prop} &= -K_{fprop} |n| n \\ M_{prop} &= -K_{mprop} |n| n \end{aligned} \quad (3)$$

$$\begin{aligned} L_{fin} &= K_L |\delta_{fin}| \delta_{fin} v_e^2 \\ M_{fin} &= K_M |\delta_{fin}| \delta_{fin} v_e^2 \end{aligned} \quad (4)$$

A more accurate thruster model can be found in [36] with the inclusion of the motor model and fluid dynamics. However, in this study, the more conservative model from [33] is used.

### 3 LOS guidance law mathematical background

This section describes the mathematical background of the 3D guidance law employ in the present study for control of an underactuated AUV. In this study, it was decided to employ the LOS proposed in [14] as it is an extension of the work of [33].

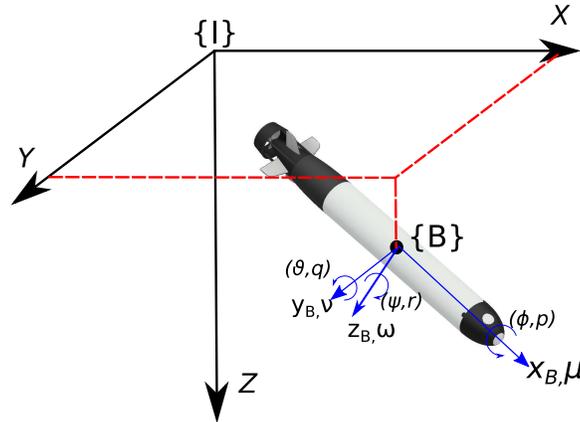


Figure 2: LOS reference frames

If the vehicle kinematic is represented by its spatial position  $p(t) \triangleq [x(t), y(t), z(t)]^T$  and its velocity is represented by  $v(t) \triangleq \dot{p}(t) \in \mathbb{R}^3$  state is related to the  $\{I\}$  frame. Also, the speed is represented by  $U(t) \triangleq |v(t)| = \sqrt{\dot{x}(t)^2 + \dot{y}(t)^2 + \dot{z}(t)^2} > 0$ . The steering of the underactuated AUV is characterized by the azimuth angle  $\chi$  and elevation  $v$  (Fig. 2).

$$\begin{cases} \chi = \text{atan2}(\dot{y}, \dot{x}) \\ v = \arctan\left(\frac{-\dot{z}}{\sqrt{\dot{x}^2 + \dot{y}^2}}\right) \end{cases} \quad (5)$$

If its consider a continuously path parametrized by a scalar variable  $\varpi \in \mathbb{R}$ , the position of a point over the path is represented by  $p_p(\varpi) \in \mathbb{R}^3$  (Fig. 3). Similarly, the orientation of the point can be defined as

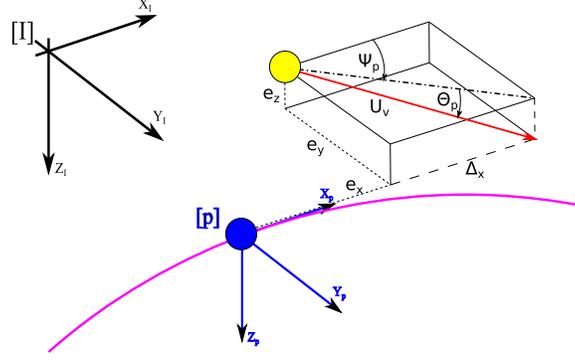


Figure 3: LOS main variables

$$\begin{cases} \psi_p = \text{atan2}(y'_p, x'_p) \\ \theta_p = \arctan\left(\frac{-z'_p}{\sqrt{x'^2_p + y'^2_p}}\right) \end{cases} \quad (6)$$

where  $x'_p = dx_p/d\varpi$ ,  $y'_p = dy_p/d\varpi$  and  $z'_p = dz_p/d\varpi$ . Hence the tracking error is expressed as:

$$\varepsilon = [x_e, y_e, z_e]^T = \mathbf{R}_F^T (p - p_p) \quad (7)$$

where  $\mathbf{R}_F^T := \mathbf{R}_z(\psi_p) \mathbf{R}_y(\theta_p)$

$$\mathbf{R}_z = \begin{bmatrix} \cos(\psi_p) & -\sin(\psi_p) & 0 \\ \sin(\psi_p) & \cos(\psi_p) & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (8)$$

$$\mathbf{R}_y = \begin{bmatrix} \cos(\theta_p) & 0 & \sin(\theta_p) \\ 0 & 1 & 0 \\ -\sin(\theta_p) & 0 & \cos(\theta_p) \end{bmatrix} \quad (9)$$

If the following Lyapunov control function positive definite is

$$V_\varepsilon = \frac{1}{2} \varepsilon^T \varepsilon \quad (10)$$

The derivative of Eq. (10) can be written as

$$\begin{aligned} \dot{V}_\varepsilon &= x_e (U_d \cos(\psi_r) \cos(\theta_r) - U_p) + \\ &+ y_e U_d \sin(\psi_r) \cos(\theta_r) - z_e U_d \sin(\theta_r) \end{aligned} \quad (11)$$

where  $U_d$  is the desired composite speed of the AUV. The auxiliary control input of the virtual point  $P_p$  is chosen as:

$$U_p = U_d \cos(\psi_r) \cos(\theta_r) + k_x x_e \quad (12)$$

where the steering angles are:

$$\begin{cases} \psi_r = \arctan\left(\frac{-k_y y_e}{\Delta_y}\right) \\ \theta_r = \arctan\left(\frac{k_z z_e}{\Delta_z}\right) \end{cases} \quad (13)$$

If the guidance variables  $\Delta_y$  and  $\Delta_z > 0$ , the control gains  $k_x, k_y, k_z$  are positive constants. If Eq. (12) and Eq. (13) are substituted into Eq. (11) and considering the relationship among inertial frame  $\{I\}$ , flow frame  $\{W\}$  and path frame  $\{F\}$ , the desired azimuth angle  $\nu_d$  and elevation angle  $\chi_d$  can be written as [37]:

$$\begin{aligned} \nu_d &= \arcsin(\sin \theta_p \cos \psi_r \cos \theta_r + \cos \theta_p \cos \theta_r) \\ \chi_d &= \text{atan2}(\chi_{d_y}, \chi_{d_x}) \end{aligned} \quad (14)$$

where

$$\begin{aligned} \chi_{d_y} &= \cos \psi_p \sin \psi_r \cos \theta_r - \sin \psi_p \sin \theta_p \sin \theta_r \\ &+ \sin \psi_p \cos \theta_p \cos \psi_p \cos \theta_r \end{aligned} \quad (15)$$

$$\begin{aligned} \chi_{d_x} = & -\sin \psi_p \sin \psi_r \cos \theta_r - \cos \psi_p \sin \theta_p \sin \theta_r \\ & + \cos \psi_p \cos \theta_p \cos \psi_p \cos \theta_r \end{aligned} \quad (16)$$

In order to transform the path following to path tracking the path was defined over time together with Eq. (14) and Eq. (12). This was to produce not only the desire angles but also the required speed at each time instance. The steering error vector can be expressed as  $[e_\mu, e_v, e_\chi]$  where  $e_\mu = \mu_d - \mu_v$ ,  $e_v = v_d - v_v$  and  $e_\chi = \chi_d - \chi_v$ .

## 4 Probabilistic Inference for Learning Control (PILCO)

PILCO algorithm [31, 38] employs GPs as the base for policy search. A GP can be defined by a mean function  $m(\cdot)$  and a positive definite covariance function  $k(\cdot, \cdot)$  commonly known as kernel. Usually a prior mean function  $m \equiv 0$  and an exponentiated quadratic kernel (Eq. (17)) are employ. This kernel only has two parameters to learn,  $l$  that determines the length of the 'wiggles' in the function and  $\sigma^2$  which determines the average distance of the function away from its mean [39].

$$k_{SE}(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right) \quad (17)$$

Given  $n$  training inputs  $X = [x_1, \dots, x_n]$  and corresponding training targets  $Y = [y_1, \dots, y_n]$ , the posterior GP hyperparameters  $l$  and  $\sigma^2$  are learned by evidence maximization [40]. The posterior predictive distribution  $p(f_* | x_*)$  of the function value  $f_* = f(x_*)$  for a test input  $x_*$  is Gaussian with mean and variance

$$\begin{aligned} m_f(x_*) &= \mathbf{k}_*^T (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y} \\ \sigma_f^2(x_*) &= k_{**} - \mathbf{k}_*^T (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_* + \sigma_\epsilon^2 \end{aligned} \quad (18)$$

```
[1] init: Sample controller parameters  $\theta \sim \mathcal{N}(0, I)$ ;
[2] Apply random control signals and record data.;
[3] repeat
[4]   Learn probabilistic (GP) dynamics model;
[5]   repeat
[6]     Approximate inference for policy evaluation;
[7]     Gradient-based policy improvement;
[8]     Update parameters  $\theta$ ;
[9]   until Convergence;
[10]  return  $\theta^*$ ;
[11]  Set  $\pi^* \leftarrow \pi(\theta^*)$ ;
[11]  Apply  $\pi^*$  to system and record data;
until Task Learned;
```

**Algorithm 1:** PILCO Algorithm

To find a policy  $\pi^*$  which converges to the desired target, PILCO builds a probabilistic GP dynamics model. This model will be the base for the deterministic approximate inference and policy evaluation, followed by the analytic computation of the policy gradients  $\partial J^\pi(\theta)/\partial\theta$  for policy improvement. The policy  $\pi$  is improved based on the gradient information  $\partial J^\pi(\theta)/\partial\theta$ .

## 5 Simulation Setup

The underwater vehicle should be able to switch between different controllers in a single mission depending on the task to be solved at a specific time. The possible variants of controllers that an AUV can use in a mission are: follow bottom, depth control, follow pipe, go-to-point, path tracking, path following, etc. A series of different simulations were designed to evaluate the capability of PILCO to control an underactuated AUV within the three different evaluation scenarios: way-point-tracking, Depth Control and Path-Tracking.

### 5.1 Vehicle Model

The vehicle model employed in this research is derived from semi-empirical calculation of the coefficients for the vehicle MULLAYA. MULLAYA is an underactuated AUV designed by the Defence Science and Technology Group

as a research platform. The vehicle is controlled by a single propeller, a pair of elevator fins and a pair of rudder fins. General specifications of the vehicle are given in Table 1. The coefficients were calculated with the same technique of [34] and the obtained coefficients are presented in Table 2. As a engineering research platform, the vehicle will undergo multiple transformations over time in shape and internal engineering. This constant change requires the needs of constant update of the controller for the vehicle. Possible future modification can include vectorized propulsion systems and buoyancy controllers.

Table 1: MULLAYA AUV particulars.

| Property | Value   | Unit       |
|----------|---------|------------|
| Length   | 1.56    | <b>m</b>   |
| Diameter | 150     | <b>mm</b>  |
| Max RPM  | 3000    | <b>RPM</b> |
| Weight   | 239.364 | <b>N</b>   |
| Buoyancy | 246.2   | <b>N</b>   |

Table 2: MULLAYA AUV coefficients employed in simulations.

| Coeff. | Result | Coeff. | Result    | Coeff. | Result |
|--------|--------|--------|-----------|--------|--------|
| Xuu    | -2.8   | Zrp    | 0.68      | Nuv    | -30.88 |
| Xwq    | -29.6  | Yuudr  | 12.12     | Npq    | -5.06  |
| Xqq    | -0.68  | Nuudr  | -7.51     | Ixx    | 0.083  |
| Xvr    | 29.6   | Zuuds  | -12.12    | Iyy    | 3.08   |
| Xrr    | -4.95  | Kpp    | -1.30E-01 | Izz    | 3.08   |
| Yvv    | -95.37 | Kpdot  | 1.09E-02  | Nwp    | 0.66   |
| Yrr    | -2.45  | Mww    | 6.96      | Nur    | -5.31  |
| Yuv    | -32.9  | Mqq    | -135.04   | Xudot  | -0.51  |
| Ywp    | 29.64  | Mrp    | 5.1       | Yvdot  | -29.64 |
| Yur    | 7      | Muq    | -5.34     | Nvdot  | 0.66   |
| Ypq    | 0.68   | Muw    | 27.16     | Mwdot  | 0.66   |
| Zww    | -95.37 | Mwdot  | -0.68     | Mqdot  | -5.05  |
| Zqq    | 2.45   | Mvp    | -0.68     | Zqdot  | -4.94  |
| Zuw    | -32.9  | Muuds  | -7.738    | Zwdot  | -29.64 |
| Zuq    | -7     | Nvv    | 6.96      | Yrdot  | 4.95   |
| Zvp    | -29.64 | Nrr    | -135.03   | Nrdot  | 5.05   |

## 5.2 Waypoint Tracking

The first type of controller evaluated for the underwater vehicle is a waypoint tracking controller. These types of controllers can be employed to transfer the vehicles between location were more specific controllers are employed or if can be work with the results of a path planning law that convert a desire path in waypoints. In this case, it is desired to control the vehicle velocity, azimuth and elevation such that the vehicle moves towards a specific location. The methodology used in this simulation is similar to the one employed in [41] but extended to 3D. If  $\eta_v$  is the vehicle state vector  $[X_v, Y_v, Z_v, \phi_v, \theta_v, \psi_v]$  that can be divided in position vector  $\mathbf{X}_v$  and orientation vector  $\theta_v$  on earth frame and the target point is expressed as the vector  $\mathbf{X}_T = [X_d, Y_d, Z_d]$  the angles of the vector between the current position and the desired position can be expressed as

$$\begin{aligned} \psi_d &= \tan^{-1} \left( \frac{Y_d - Y_v}{X_d - X_v} \right) \\ \theta_d &= \tan^{-1} \left( \frac{\sqrt{(Y_d - Y_v)^2 + (X_d - X_v)^2}}{Z_d - Z_v} \right) \end{aligned} \quad (19)$$

and the vector of angle errors can be expressed as the difference between the vehicle orientation and the desire orientation, i.e  $\theta_e = [e_\psi, e_\theta, e_u]$ , where  $e_\psi = \psi_d - \psi_v$  and  $e_\theta = \theta_d - \theta_v$ . The target of the policy will be to minimize the error angles to zero and the surge speed error to zero. In the proposed simulation a surge speed of  $1.2m/s$  was set as the desired speed with an initial position of  $[0, 0, 0]$ , an orientation of  $[0, 0, \pi/4]$  and an initial surge speed of  $0.5m/s$ . The first 12 seconds of real model simulation was employed to learn the policy with a target point with coordinates

[40, 40, 10]. A constraint to limit the vehicle turn  $< 180^\circ$  was applied to the process of policy testing over the vehicle model. A second simulation with two target points [30, 30, 10] and [90, 100, 10] is employed to check the viability of the policy. A total of 1000 sparse points was located as the limiter from which a sparse model will be employ. The policy and simulation were executed at  $5Hz$ . A noise with variance of  $\sigma^2 = 0.005$  was employed in the simulation and a variance of  $\sigma^2 = 0.2$  was employed for the start position of the vehicle. When the vehicle arrives at 3 meters of the objective the vehicle will be given the secondary target as the new objective. If the vehicle arrives at 1 meter of the final target the simulation will stop.

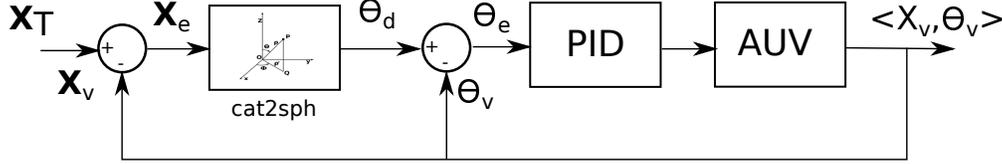


Figure 4: Control system block for waypoint tracking with PID controller.

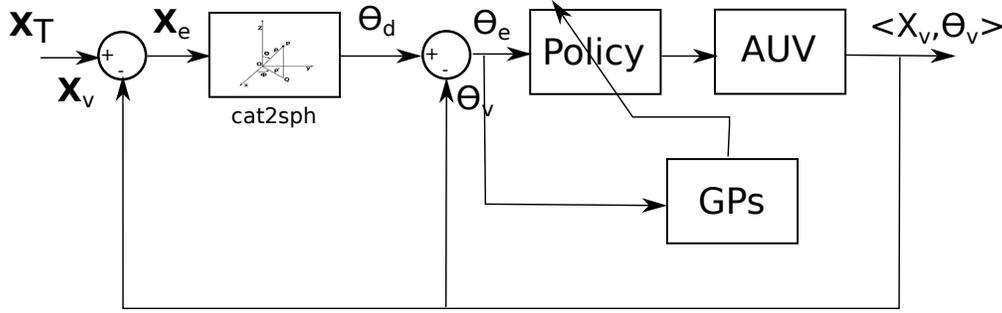


Figure 5: Control system blocks for waypoint tracking with RL policy.

### 5.3 Depth Control

A secondary scenario is explored to evaluate the controller’s capability to keep a specific depth while traveling between location. This is an extension to the previous search of a policy, i.e. waypoint. In this scenario, we want the vehicle to arrive at a specific depth before going to the desired location. The policy will minimize the vector of error to zero with the difference that the vector of errors is  $\theta_e = [e_\psi, e_\theta, e_u, e_z]$ , where  $e_z = z_d - z_v$  similar to the waypoint tracking. The vehicle starts from the same initial position of the waypoint tracking scenario [0, 0, 0]. the proposed training target was the vector [40, 40, 5] and the selected test target points are [30, 30, 5] and [90, 100, 5]. A total of 1500 sparse points was located as the limiter from which a sparse model will be employed. The policy and simulation were executed at  $5Hz$ . A noise of  $\sigma^2 = 0.005$  was employed for the simulation and a  $\sigma^2 = 0.2$  was employed for the start position of the vehicle. When the vehicle arrives at 3 meters of the objective the vehicle will be given the secondary target as the new objective. I f the vehicle arrives at 1 meter of the final target the simulation will stop.

### 5.4 Path Tracking

With the aim to test the capability of PILCO for path tracking, the scenario consists of a single policy to control propeller force, elevator force and rudder force of the vehicle. The decision to learn the force and not to control direct the RPM and angle of the fins was to be able to compare the policy to a standard controller from the literature. In the design of the authors path tracking simulation, the equation presented in section Section 3. Fig. 7 shows the block diagram of the LOS-PILCO control implementation whereby the policy will evolve based on the learned GPs. The target of all learned policies is to reduce the vector of errors  $[e_\mu, e_v, e_\chi]$  to zero. A LOS-PID was used as a performance comparison. The LOS-PID (Fig. 6)with the exact coefficient of the model was coded in the same ways as [14] with the inclusion of measurement noise. The initial position of the vehicle was established as [60, 3, 1] with an initial orientation [0, 0,  $3\pi/4$ ]. A total of 1500 sparse points was located as the limiter from which a sparse model will be employed. A noise with  $\sigma^2 = 0.001$  was employed throughout the simulation and a random variation of the start point with a variance of  $\sigma^2 = 0.2$  was employed. Higher values of  $\sigma$  were not possible for the LOS-PID to be comparable

as it was not able to overcome higher values of noise. Both the PID and PILCO policy were executed at a frequency of 10 Hz. A helix path (Fig. 8) was parametrized as is show in Eq. (20).

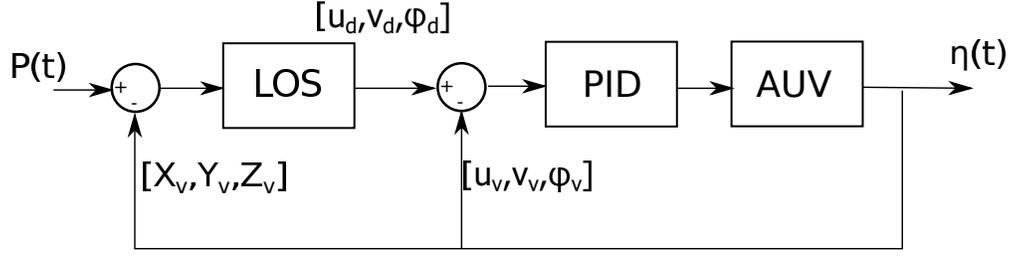


Figure 6: Control systems block for LOS-PID

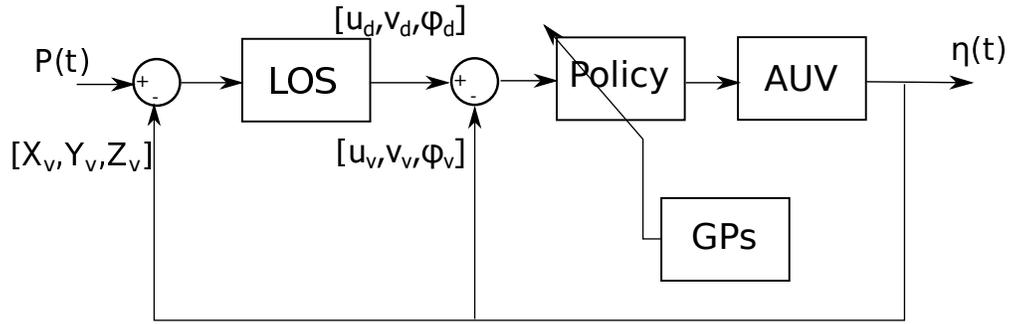


Figure 7: Control systems block for LOS-PILCO

$$\begin{aligned}
 X_{helix} &= 60 * \cos(0.02618 * w_{ramp}(t)) \\
 Y_{helix} &= 60 * \sin(0.02618 * w_{ramp}(t)) \\
 Z_{helix} &= 2 + 2 * w_{ramp}(t)/200
 \end{aligned} \tag{20}$$

where  $w_{ramp}(t)$  is a function over time with a slope  $m$ . The control learning of the vehicle was done over the first 20 seconds at a frequency of 10 Hz.

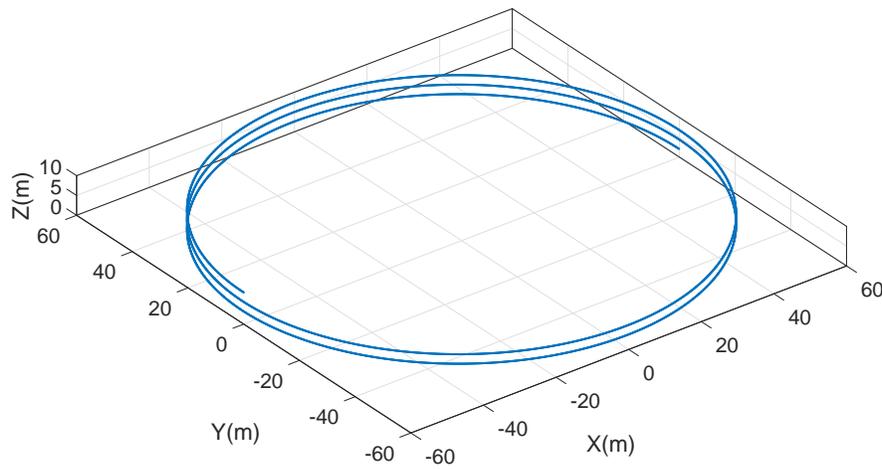


Figure 8: Spiral Path to be follow

## 6 Results

### 6.1 Waypoint Tracking

As evident in Fig. 11, the learning of a policy for waypoint tracking of an AUV is possible with the application of PILCO. However, reinforced learning does not require the calibration of parameters, but rather, requires the tuning of the parameters of the  $Q$  matrix from the cost function. For the cost function design, the surge speed was given a higher importance in the cost function than the other error vectors. This escalation of each target in the learning policy was needed as an underactuated AUV needs to get to a specific speed to be able to dive and navigate with a more linear model. An example of the data employed to create the GPs model can be seen in Fig. 9.

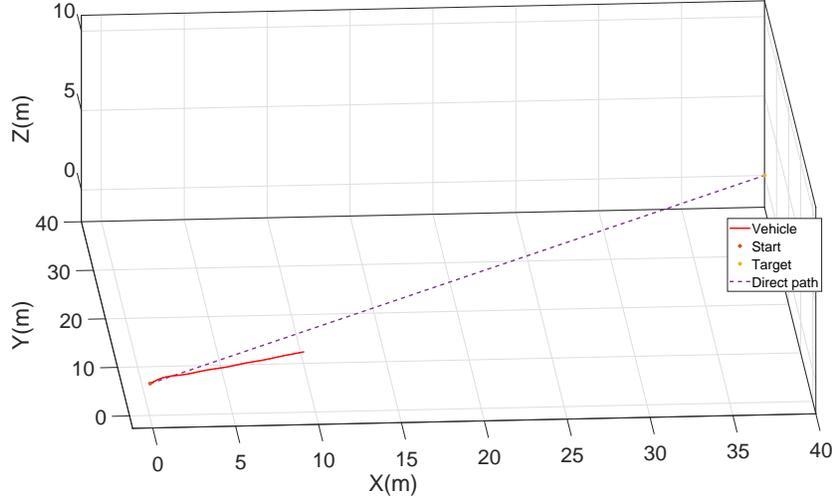


Figure 9: Real model waypoint tracking data example for GPs learning.

The authors simulations showed that the required number of episodes over the platform was under 25 to obtain a usable policy to control the vehicle and be able to arrive at all three targets. Fig. 10 shows the evolution of the cost function over time for the platform. The cost function drops rapidly to a low value as the vehicle learns to control the surge speed and from there, how to control its orientation. The decision to do this is based in that the control surfaces as fin require a minimum speed to be useful.

The result from the final policy selected is shown in Fig. 11; plotting against the results from the PID controller for the same scenario. Both PID and PILCO arrive near to both test point. However, PILCO shows a more direct route taken towards the targets. The PID controller took 1201 cycles to arrive equivalent to 240.2 seconds and the PILCO policy took 569 cycles equivalent to 113.8 seconds. PILCO has an advantage over a simple PID as PILCO has learned the policy over a noisy platform (or environment) and the PID doesn't have any component to compensate for the noise in the measurement. For PID to compensate for the noise an observer will have to be designed and implemented. However, this will increase the complexity of the controller and the need for more design and calibration time for the PID controller. The results show that the waypoint tracking with a PILCO controller can be a viable option as the tuning of the controller coefficients are practically zero. The robustness of the PILCO policy is higher as the platform will start the learning from different angles and can overcome noise in the measurement. The policy is also updated in case of failure. A down point of PILCO in the simulated scenarios is that the controller tries to always have the same behaviour, e.g. if the vehicle dives a little before directing to the first point after the second point the policy will try to repeat the same behaviour.

### 6.2 Depth control

The simulation of simultaneous waypoint tracking, and depth control of underwater vehicles shows that PILCO can learn a complex policy to control the vehicle. An example of the training data used for learning the vehicle GPs model can be seen in Fig. 12. The quality of the learning of the GPs model is directly related to the quality of the learning of the policy. After 20 episodes over the platform a good quality GPs model has been learned such that a policy can be learnt. The evolution of the cost function over time can be seen in Fig. 13. In similar way to the standard

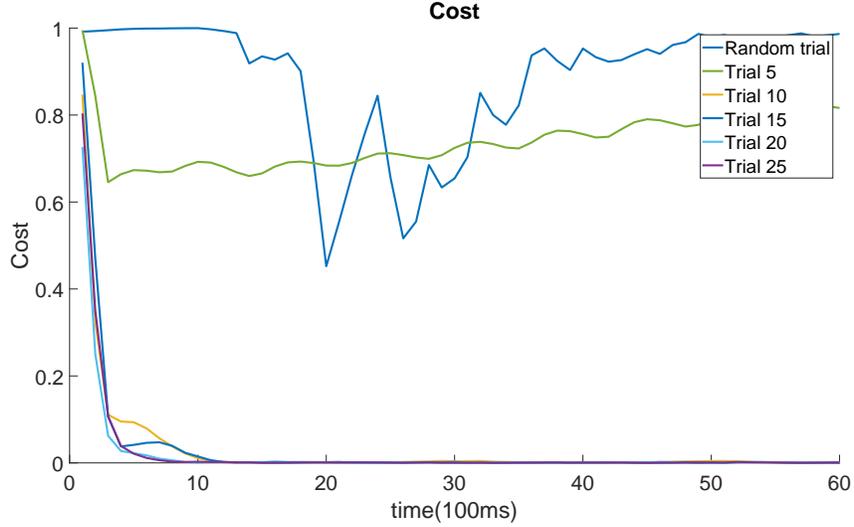


Figure 10: Waypoint tracking cost function evolution.

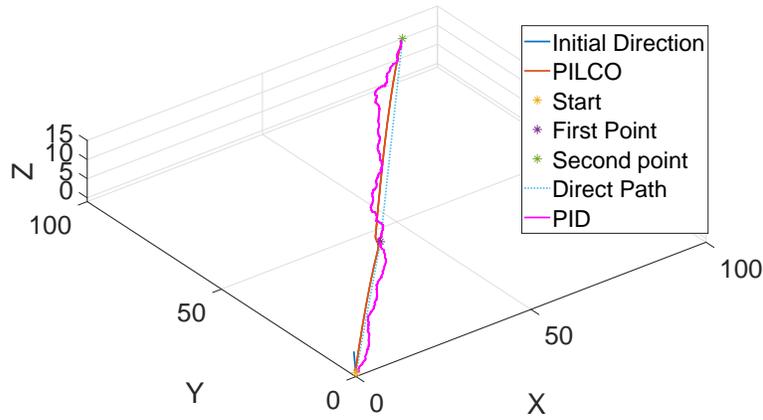


Figure 11: 3D view Waypoint tracking comparative result between PILCO and PID Results [m].

waypoint tracking, the cost evolves very fast to a low value and then small changes are done in the policy learning. This small change is produced by the scaling of each of the targets to be obtained.

Fig. 14 shows the results of policy 22. The policy is capable of keeping the depth at near to 5 meters and at the same time keep all four errors at near zero value. Usually, the task described here will be done with two controllers a depth controller and a azimuth controller. As we are applying a single controller with two targets(elevation, depth) that at the start go on different paths. The  $[X, Y, Z]$  results from the simulation are presented in Fig. 15. The policy produces a different behaviour than normal controllers, if by operator error a target was setup with a different depth than the desire depth the learned policy will try to push the vehicle to the target disregarding the depth by a small quantity such the target can be completed this can be observed in Fig. 16 . Three targets were setups as  $[30, 30, 5]$ ,  $[70, 70, 3]$ , and  $[90, 100, 5]$  the second target is outside of the desire depth of the policy.

### 6.3 Path Tracking

The training evolution of the GPs model in the direction of a better policy allows PILCO to search for a better policy after each iteration with the real model (Fig. 18). Fig. 17 presents an example of the data used for training of the GPs model that is later used for policy search. The reinforced learning simulation to learn a LOS controller with PILCO shows that the reinforced learning algorithm can learn to follow the desired path based on the desire angles produced by a LOS law. The policy can perform better than a LOS-PID controller that is without noise compensation algorithms.

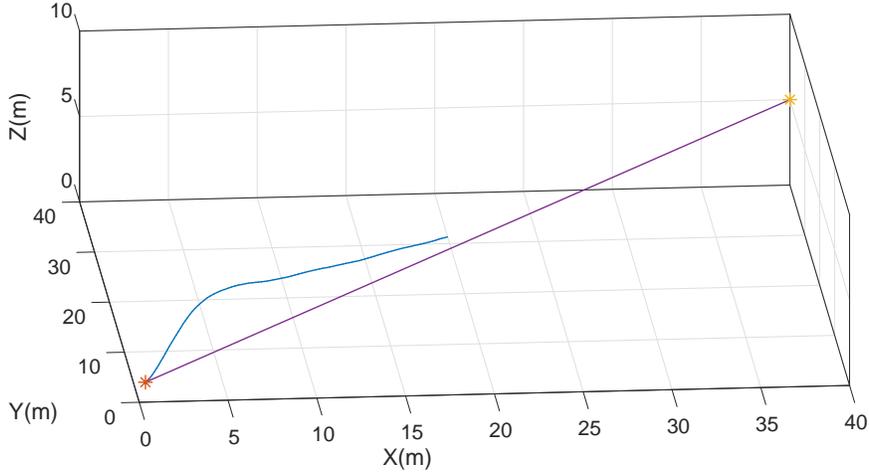


Figure 12: Real model waypoint tracking and depth control data example for GPs learning.

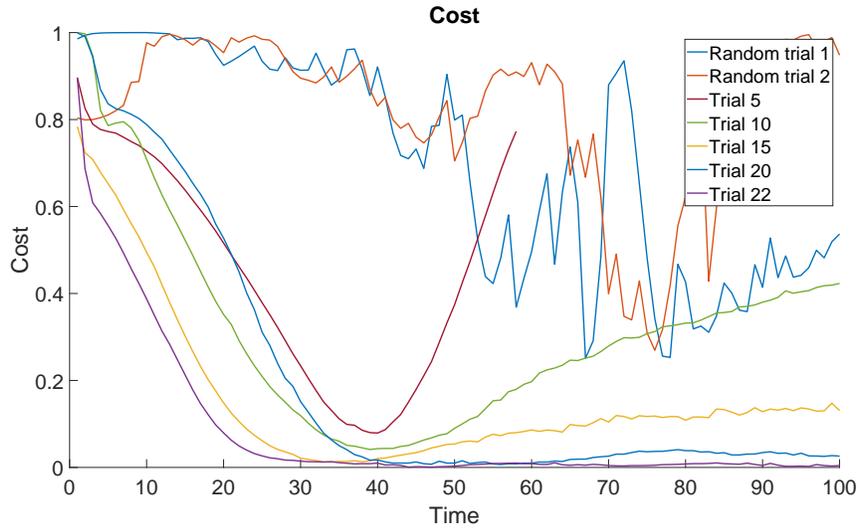


Figure 13: Waypoint tracking and depth control cost function evolution.

The results from the PID to follow the LOS law can be seen in Fig. 19, without noise the PID can stabilize himself very fast but with minimum noise in the depth sensor the  $\chi$  angle cannot be completely stable and the controller will fluctuate. In the same way, PILCO controller (Fig. 20) reacts to the noise of the measurement but can follow the desired path in a better way that the LOS-PID controller implemented.

The measurement of the RMSE between the desire vehicle position and the vehicle position for both LOS-PID and LOS-PILCO is presented in Table 3. LOS-PILCO is shown to out-perform the PID controller in the reduction of error between the LOS law and the vehicle position. Not only in the accuracy of placing the vehicle in the correct position but in the speed of deployment of a controller to follow LOS will LOS-PILCO will require less tuning with field tests. Another advantage of PILCO is the absence of knowledge from the vehicle coefficients and not need to use other vehicles variables a surge, heave and sway speeds.

Fig. 21 and Fig. 22 present the comparative plots of position in  $X, Y, Z$  and the 3D path taken by the vehicle with LOS-PID controller and LOS-PILCO policy. Both controllers shows their ability to direct the AUV to the desired path, with the PILCO policy showing a better performance after episode 20.

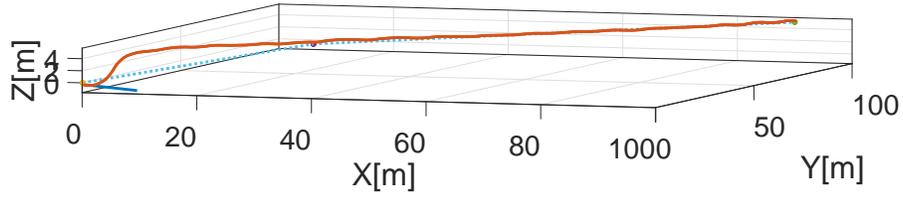


Figure 14: 3D view Waypoint tracking and depth control PILCO results.

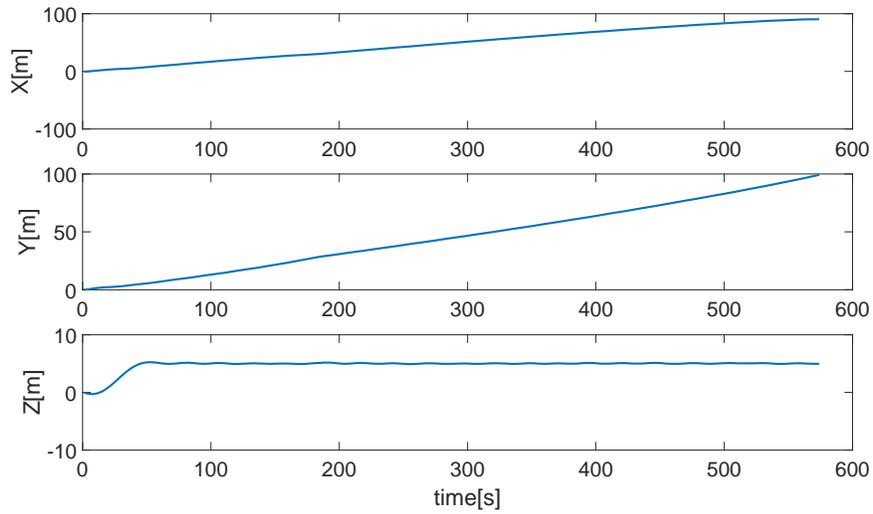


Figure 15:  $[X, Y, Z]$  Waypoint tracking and depth control PILCO results.

Table 3: RMSE results from LOS-PID and LOS-PILCO.

| RMSE      | Value |
|-----------|-------|
| LOS-PID   | 1.214 |
| LOS-PILCO | 0.89  |

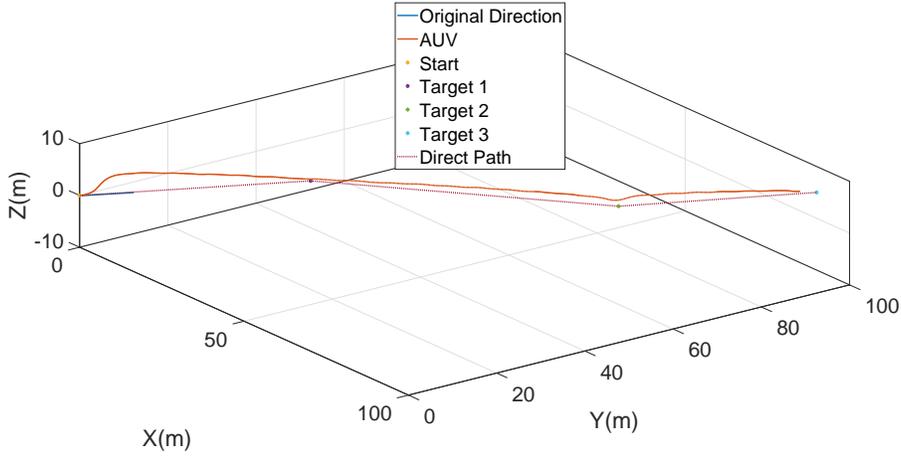


Figure 16: 3D Waypoint tracking and depth control PILCO results with a target outside of the control depth.

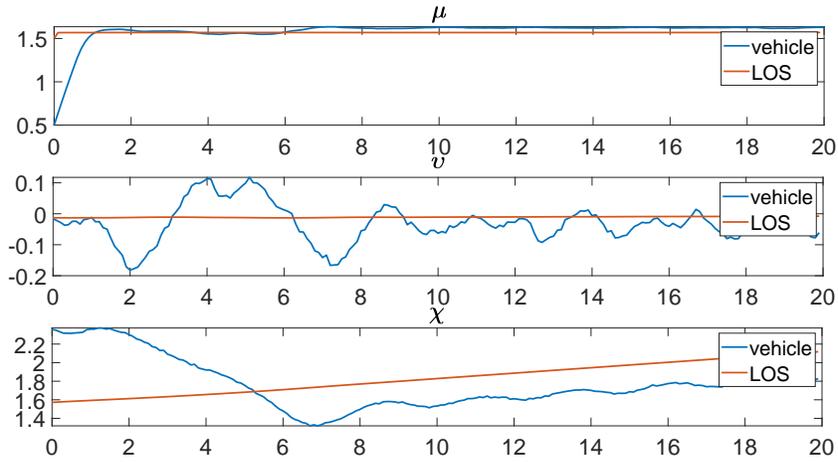


Figure 17: Sample of data used for learning of GPs model employed.

## 7 Conclusions

This paper investigates the applicability of PILCO algorithms and its requirements to learn policies to control under-actuated AUVs. Three sets of simulations were designed to evaluate the capability of PILCO for waypoint tracking, depth control and path tracking. The simulations shown that a simple waypoint tracking control can be learnt in a small quantity of experiments over the real vehicle, the performance of the learnt policy was compared with a PID controller which is over-performed by the policy as the policy obtained is learned over the platform and the non-parametric model includes noise.

In a similar way, a depth control policy was learned by mixing the waypoint tracking objective with the depth objective. However, in this research the learning of a policy to only control depth was not successful as the GPs model to allow the learning of a policy requires more information. Nevertheless, in the simulation of depth control it had been shown that a policy of depth control can be obtained by the learning of simultaneous waypoint tracking and depth control. The combination of objectives gives us a more intuitive behaviour, similar to what a human will do with the proposed objectives.

In the case of the proposed LOS-PILCO methodology, it has been shown that PILCO is a viable option to learn a policy to minimize the error between a LOS law and the vehicle position. The PILCO policy is shown to perform

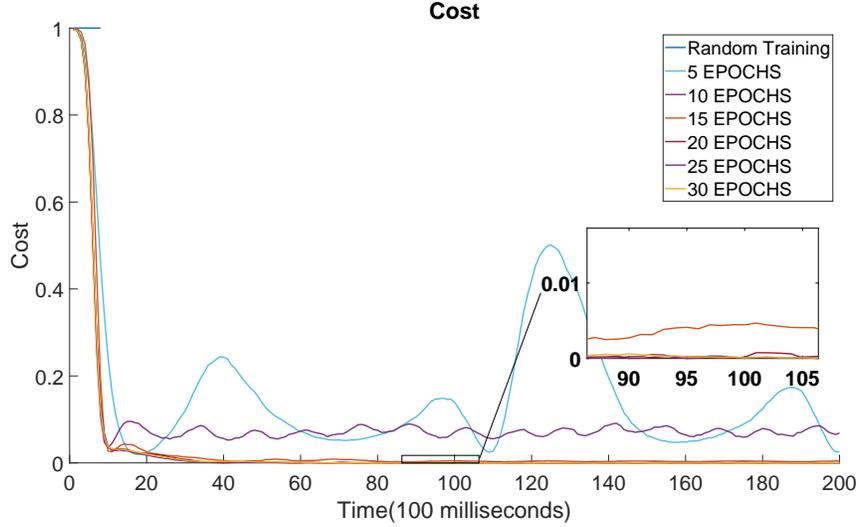


Figure 18: Cost evolution for each policy test over real model.

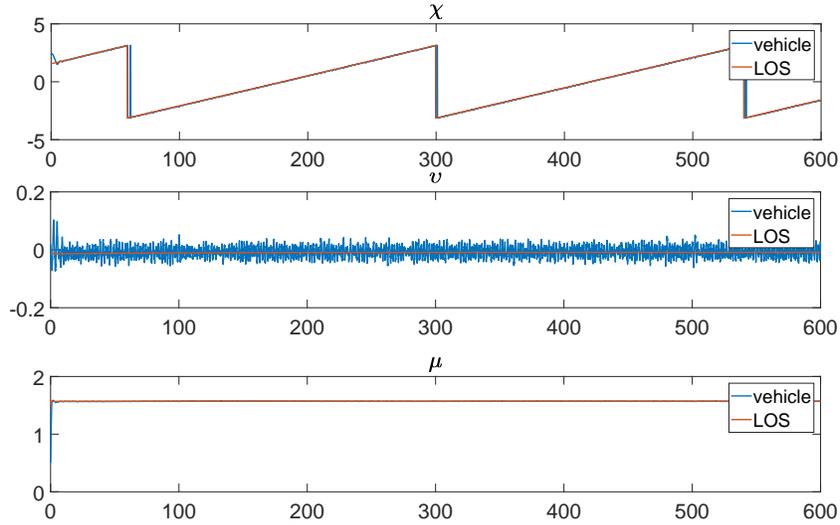


Figure 19: Desire LOS angles and speed and vehicle LOS angles and speed produced by PID controller.

equally and sometimes better than a PID. The RMSE shows that PILCO can obtain better performance and the long period of simulation shows that the learned policy can constantly minimize the error to the desired value.

PILCO algorithm has shown that is applicable to an underactuated AUV. In the simulations were consider limits that are a requirement for safety of the vehicle as maximum depth and maximum angles. The limits imposed to the platform do not limited the capability of PILCO to learn a viable policy. In the design of the cost function ,the vehicle shape and type of actuators force the selection of values of the cost function such that forward speed has to represent a higher cost than the specific vehicle angles or position. Model based reinforced learning for underactuated AUV shows to be a solution motion control of underwater vehicles and can be a solution to control bioinspired vehicles which are more dynamically complex to describe with a mathematical model.

## References

- [1] Zaopeng Dong, Lei Wan, Yueming Li, Tao Liu, Jiayuan Zhuang, and Guocheng Zhang. Point stabilization for an underactuated auv in the presence of ocean currents. *International Journal of Advanced Robotic Systems*,

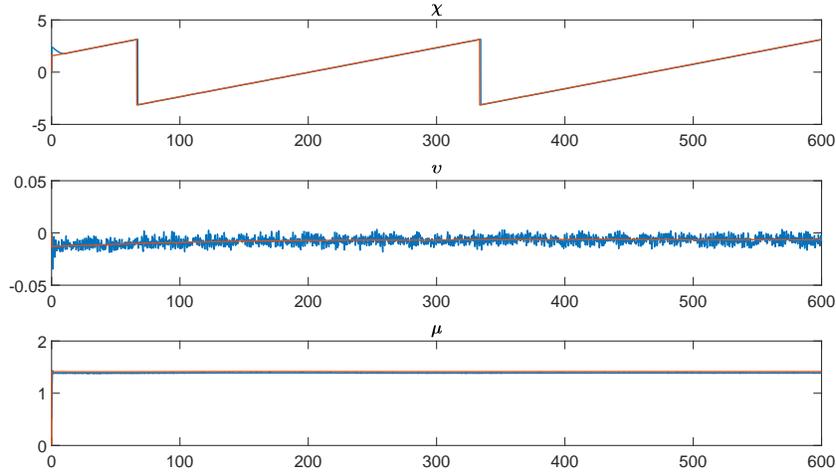


Figure 20: Desire LOS angles and speed and vehicle LOS angles and speed produced by PILCO policy.

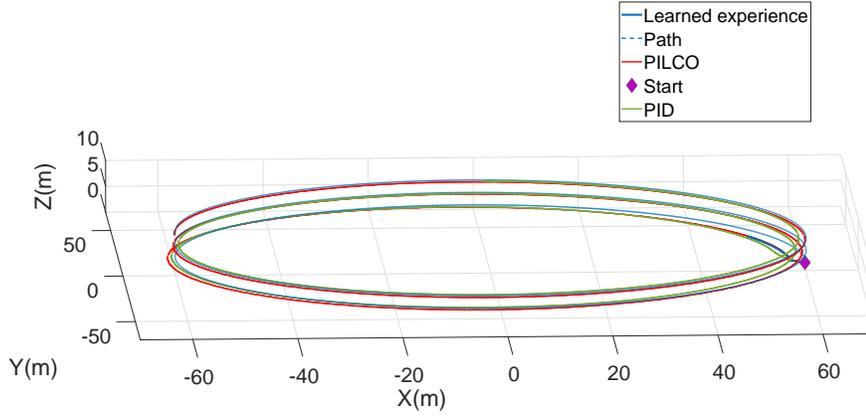


Figure 21: 3D comparison of vehicle controlled with PID and PILCO controller.

- 12(7):100, 2015.
- [2] F. Alonge, F. D’Ippolito, and F. M. Raimondi. Trajectory tracking of underactuated underwater vehicles. In *Proceedings of the 40th IEEE Conference on Decision and Control (Cat. No.01CH37228)*, volume 5, pages 4421–4426 vol.5, Dec 2001.
  - [3] M. Breivik and T. I. Fossen. Guidance-based path following for autonomous underwater vehicles. *Oceans 2005, Vols 1-3*, pages 2807–2814, 2005.
  - [4] X. Xiang, L. Lapierre, C. Liu, and B. Jouvencel. Path tracking: Combined path following and trajectory tracking for autonomous underwater vehicles. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3558–3563, Sept 2011.
  - [5] Ahmad Forouzantabar, Babak Gholami, and Mohammad Azadi. Adaptive neural network control of autonomous underwater vehicles. *World Academy of Science, Engineering and Technology*, 6(7):304–309, 2012.
  - [6] K. D. Do, J. Pan, and Z. P. Jiang. Robust and adaptive path following for underactuated autonomous underwater vehicles. *Ocean Engineering*, 31(16):1967–1997, 2004.

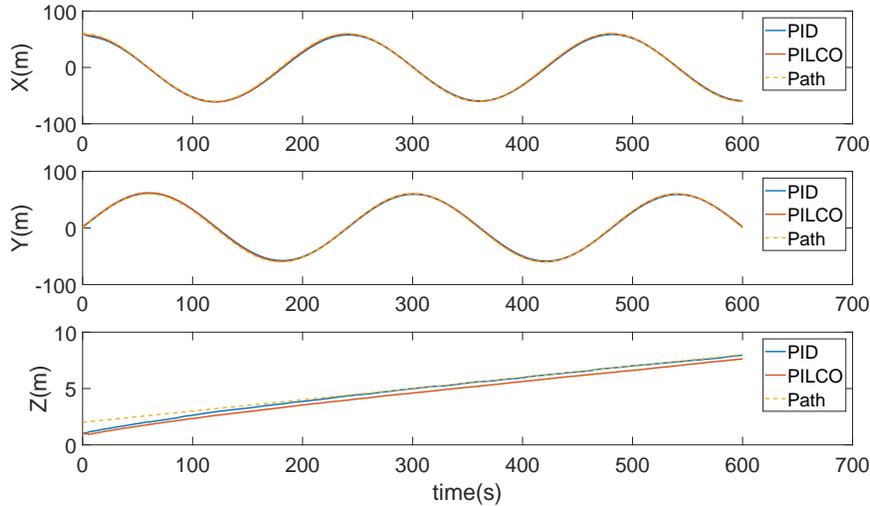


Figure 22: X,Y,Z comparison of vehicle controlled with PID and PILCO policy.

- [7] Xianbo Xiang, Lionel Lapierre, and Bruno Jouvencel. Smooth transition of auv motion control: From fully-actuated to under-actuated configuration. *Robotics and Autonomous Systems*, 67:14 – 22, 2015. Advances in Autonomous Underwater Robotics.
- [8] Chris Roper. Using the gavia auv system to locate and document munitions dumped at sea. Online, oct 2007.
- [9] Sigurd Andreas Holsen. Dune: Unified navigation environment for the remus 100 auv-implementation, simulator development, and field experiments. Master’s thesis, NTNU, 2015.
- [10] R. Rout and B. Subudhi. Narmax self-tuning controller for line-of-sight-based waypoint tracking for an autonomous underwater vehicle. *IEEE Transactions on Control Systems Technology*, 25(4):1529–1536, July 2017.
- [11] Mansour Ataei and Aghil Yousefi-Koma. Three-dimensional optimal path planning for waypoint guidance of an autonomous underwater vehicle. *Robotics and Autonomous Systems*, 67:23–32, 2015.
- [12] S Saravanakumar and T Asokan. Waypoint guidance based planar path following and obstacle avoidance of autonomous underwater vehicle. In *ICINCO (2)*, pages 191–198, 2011.
- [13] C. Yu, X. Xiang, and J. Dai. 3d path following for under-actuated auv via nonlinear fuzzy controller. In *OCEANS 2016 - Shanghai*, pages 1–7, 2016.
- [14] Xianbo Xiang, Caoyang Yu, and Qin Zhang. Robust fuzzy 3d path following for autonomous underwater vehicle subject to uncertainties. *Computers and Operations Research*, 84:165 – 177, 2017.
- [15] S. Wang, Y. Shen, Q. Sha, G. Li, J. Jiang, J. Wan, T. Yan, and B. He. Nonlinear path following of autonomous underwater vehicle considering uncertainty. In *2017 IEEE Underwater Technology (UT)*, pages 1–4, Feb 2017.
- [16] W. Caharija, K. Y. Pettersen, J. T. Gravdahl, and E. Borhaug. Path following of underactuated autonomous underwater vehicles in the presence of ocean currents. *2012 Ieee 51st Annual Conference on Decision and Control (Cdc)*, pages 528–535, 2012.
- [17] Xiao Yang, Yue Shen, Kaihong Wang, Qixin Sha, Bo He, and Tianhong Yan. Path following for an autonomous underwater vehicle using gp-los. In *OCEANS 2016-Shanghai*, pages 1–5. IEEE, 2016.
- [18] Filoktimon Repoulias and Evangelos Papadopoulos. Planar trajectory planning and tracking control design for underactuated auvs. *Ocean Engineering*, 34(11):1650 – 1667, 2007.
- [19] Xiao Liang, Yuan You, LF Su, Wei Li, and Jundong Zhang. Path following control for underactuated auv based on feedback gain backstepping. *Technical Gazette*, 22(4):829–835, 2015.
- [20] Jian Gao, Weisheng Yan, Ningning Zhao, and Demin Xu. Global path following control for unmanned underwater vehicles. In *Control Conference (CCC), 2010 29th Chinese*, pages 3188–3192. IEEE, 2010.
- [21] Xiao Liang, Lei Wan, James I.R. Blake, R. Ajit Sheno, and Nicholas Townsend. Path following of an underactuated auv based on fuzzy backstepping sliding mode control. *International Journal of Advanced Robotic Systems*, 13(3):122, 2016.

- [22] Z. Chu and D. Zhu. 3d path-following control for autonomous underwater vehicle based on adaptive backstepping sliding mode. In *Information and Automation, 2015 IEEE International Conference on*, pages 1143–1147, 2015.
- [23] Xinqian Bian, Jiajia Zhou, Zheping Yan, and Heming Jia. Adaptive neural network control system of path following for auvs. In *2012 Proceedings of IEEE Southeastcon*, pages 1–5, March 2012.
- [24] J. Wang, C. Wang, Y. Wei, and C. Zhang. Three-dimensional path following of an underactuated auv based on neuro-adaptive command filtered backstepping control. *IEEE Access*, pages 1–1, 2018.
- [25] H. Kawano. Method for applying reinforcement learning to motion planning and control of under-actuated underwater vehicle in unknown non-uniform sea flow. In *2005 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 996–1002, Aug 2005.
- [26] H. Shen and C. Guo. Path-following control of underactuated ships using actor-critic reinforcement learning with mlp neural networks. In *2016 Sixth International Conference on Information Science and Technology (ICIST)*, pages 317–321, May 2016.
- [27] Runsheng Yu, Zhenyu Shi, Chaoxing Huang, Tenglong Li, and Qiongxiang Ma. Deep reinforcement learning based optimal trajectory tracking control of autonomous underwater vehicle. In *Control Conference (CCC), 2017 36th Chinese*, pages 4958–4965. IEEE, 2017.
- [28] Sigurd A Fjerdingen, Erik Kyrkjebø, and Aksel A Transeth. Auv pipeline following using reinforcement learning. In *Robotics (ISR), 2010 41st International Symposium on and 2010 6th German Conference on Robotics (ROBOTIK)*, pages 1–8. VDE, 2010.
- [29] Christopher John Cornish Hellaby Watkins. *Learning from delayed rewards*. PhD thesis, King’s College, Cambridge, 1989.
- [30] Chris Gaskett, David Wettergreen, Alexander Zelinsky, et al. Reinforcement learning applied to the control of an autonomous underwater vehicle. In *Proceedings of the Australian Conference on Robotics and Automation (AuCRA99)*, 1999.
- [31] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian processes for data-efficient learning in robotics and control. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):408–423, Feb 2015.
- [32] M. De Paula and G. G. Acosta. Trajectory tracking algorithm for autonomous vehicles using adaptive reinforcement learning. In *OCEANS 2015 - MTS/IEEE Washington*, pages 1–8, Oct 2015.
- [33] Thor I Fossen. *Guidance and control of ocean vehicles*, volume 199. Wiley New York, 1994.
- [34] Timothy Prester. *Verification of a six-degree of freedom simulation model for the REMUS autonomous underwater vehicle*. Thesis, 2001.
- [35] Morton Gertler and Grant R Hagen. Standard equations of motion for submarine simulation. Technical report, DAVID W TAYLOR NAVAL SHIP RESEARCH AND DEVELOPMENT CENTER BETHESDA MD, 1967.
- [36] Jinhyun Kim and Wan Kyun Chung. Accurate and practical thruster modeling for underwater vehicles. *Ocean Engineering*, 33(5-6):566–586, 2006.
- [37] Morten Breivik and Thor I. Fossen. Guidance laws for autonomous underwater vehicles. In Alexander V. Inzartsev, editor, *Underwater Vehicles*, chapter 4. IntechOpen, Rijeka, 2009.
- [38] Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472, 2011.
- [39] David Duvenaud. *Automatic model construction with Gaussian processes*. PhD thesis, University of Cambridge, 2014.
- [40] Carl Edward Rasmussen. Gaussian processes in machine learning. In *Advanced lectures on machine learning*, pages 63–71. Springer, 2004.
- [41] Shusheng Bi, Chuanmeng Niu, Yueri Cai, Lige Zhang, and Houxiang Zhang. A waypoint-tracking controller for a bionic autonomous underwater vehicle with two pectoral fins. *Advanced Robotics*, 28(10):673–681, 2014.

This figure "Wilmer.jpg" is available in "jpg" format from:

<http://arxiv.org/ps/1912.11584v1>

This figure "zleong.jpg" is available in "jpg" format from:

<http://arxiv.org/ps/1912.11584v1>