# A Generic Conceptual Framework for Autonomous Navigation with a case study on Wildlife Detection using UAV

Gopalakrishnan Saisubramaniam (A0178249N), Goh Chung Tat Kenrick (A0080891Y)
Master of Technology in Knowledge Engineering, Institute of Systems Science,
National University of Singapore

*Abstract*—**The growing interest in autonomous navigation is seen through the emergence of several innovations within the last decade. From classical rule-based approaches to convolutional neural networks to reinforcement learning, the method adopted to make an agent learn to act on its own has varied based on the complexity of the task at hand and the environment it must work in. Reinforcement Learning in particular, has seen a resurgence in popularity with the successful implementation of it on classical Atari game. Based on learning through experience/ trial and error which is similar to a normal human learning process, it allows an agent having imperfect information to understand the inputs given to it through the environment and determine the best possible action that should be done at that point of time. The more number of times the agent plays, the more experience it gains and becomes better at the game. Adopting the concept for real world tasks is possible, however not in a strict sense. Time and resource constraints discourage training an agent from scratch. Hence simulations are often adopted to bootstrap the training process. We propose a conceptual framework general enough to allow the agent to learn faster by defining a boundary on the search space and rewarding it to achieve maximal coverage of the objective set for the agent. As a case study, we simulate an autonomous drone in an African environment setting to fly over a given set of source and destination points and achieve the set objective of detecting wildlife from a fixed height and distance while minimizing/avoiding collision with the artifacts of the environment.**

*Index Terms*— **Autonomous navigation, reinforcement learning, computer vision, deep learning, artificial intelligence**

## I. INTRODUCTION

The motive behind autonomous movement is to enable a machine powered by artificial intelligence to navigate around on its own without external supervision, relying on its own physical and intellectual capabilities. Technical advancements have been made since the last three decades, starting with rule-based approaches, [12] Kalman filters [13], convolutional neural networks [14] and more recently using reinforcement learning techniques. [15]

.

Reinforcement learning has been around since 1998 [1] and it has been receiving increasingly amount of attention in the past 4 years with the advancement in computation power. Reinforcement learning is a part of the multi-disciplinary field namely psychology and neuroscience as the agent attempts to learn through experience and then devise new strategies to improve on its action in the environment. Computer games allows the agent to work in a virtual environment where all the variables can be controlled by the game engine. In addition, considering from the viewpoint of a computer game, there is a comparative benchmark available easily in all games. By pitting the artificial agent against Real Players or its own, we can determine if the model is working well or not. Some examples of successful Deep Reinforcement Learning done is Alpha GO where AlphaGO [3], the artificial agent is able to defeat the Grand Master using Reinforcement Learning. In 2017, Reinforcement Learning has been expanded to more complex games such as Starcraft [4]. In Starcraft, it is a unique concept as it is a multiple agent problem with vastly large search space and imperfect information. Although the project is not able to succeed in creating an artificial agent to replicate a pro-player's level of expertise, it is deemed as a prelude to the next major achievement. In 2018, OpenAI Five has successfully create 5 artificial agents to co-operate with each other and compete with a human team of 5. It is able to defeat a semi-pro team which is a huge achievement, considering the imperfect information as well as the vast search space [5].

## II. RELATED WORK

In 2015 Deep Q-Learning (DQN) [2] was published which is able to handle a variety of complex tasks normally requiring human level control. In particular, DQN is able to learn from its own experience in all the classic Atari 2600 games and achieve a comparable result compared to a professional human game tester. The work bridges the divide between high-dimensional sensory inputs and actions leading to the first artificial agent capable of learning to excel at a diverse array of challenging task.

### III. TERMINOLOGY

A. **Autonomous Agent** – a self-sufficient system independent of an external server/aid that has the capabilities to perceive sensory and visual input at state *s*, process and compute inference on the fly and take one of the predefined actions *a* and attempt to maximize its reward *r* during the course of its operation to achieve the defined objective *o* set for the episode *e*.

B. **State** – current location of the agent in the search space when receiving its input.

C. **Action** – the set of all possible moves performed by agent, predefined before the start of the algorithm.

D. **Policy** – Policy refers to the strategy that the agent will employ to determine the next action based on current state. It aims to choose the next action which provides the best rewards.

E. **Reward** – Reward refers to the short-term output given by the environment based on the last action. The agent learns to maximize the cumulative future reward which is denoted in the equation below; can be positive (good feedback) or negative(penalty)

$$R_t = r_{t+1} + r_{t+2} + r_{t+3} \ldots = \sum_{k=0} r_{t+k+1}$$

F. **Objective** – the business goal to be achieved a.k.a. the purpose of the agent to exist.

G. **Episode** – the interval from the source to the destination; comprises of several states and actions between the start and end.

H. **Value** – Refers to the expected long-term return which differs from the short-term reward R. It is based on the state that the agent finds itself while choosing different actions during the training progress. It is noted that while the state may be the same for different circumstances, Different actions may take for different iterations in order to arrive at that state. As such, the reward R will differ slightly and have some randomness in it. As such, the expectation takes into account all this randomness and comes out with a mean value.

$$V^\pi(s) = \mathbb{E}_\pi\big[R_t | s_t = s\big]$$

I. **Q-value** – Refers to the long-term return of the current state when taking action under a policy. This is similar to the variable Value. However, it tells us the expected Cumulative Reward given an action a. We aim to maximize this Q-value which gives us an indicator that the agent is being more proficient at performing the task

$$Q^\pi(s, a) = \mathbb{E}_\pi\big[R_t | s_t = s, a_t = a\big]$$

A. **Environment** – The environment represents the virtual world where the agent performs an action. The environment will take the actions and returns 2 outputs, namely the agent's reward as well as the current state in the environment.

Overall, the diagram below shows the architecture on how the whole algorithm works and how the agent learns from its action given the rewards and states which are output by the environment.
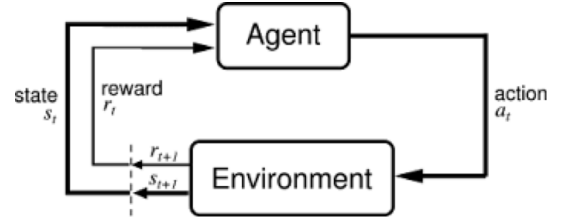


**Figure 1** Agent Environment Interaction

### IV. PROPOSED CONCEPTUAL FRAMEWORK

A. *Overview*

We propose a framework that takes into consideration the navigation, objective satisfaction and route tracing aspects of an agent.

B. *Business Objectives*

Although not an exhaustive list, some of these objectives can be covered in an autonomous navigation project:
   i. Cost – the cost assigned to travel a certain distance. The total number of steps taken by agent to reach the destination. Steps that are further away have either a lower positive or a higher negative cost and viz with the steps closer to the destination.
   ii. Constraints – A set of constraints must be enforced on the agent w.r.t the environment designated to it. The viewpoint is to effectively reduce the search space using these constraints.
   iii. Target – There must be a target for the agent to achieve, using which the agent reinforces its belief states and progresses towards the objective.

C. *Navigation Algorithm*

The algorithm suggested for navigation is DQN reinforcement learning proposed in the initial paper [2]. At each time step, it chooses an action $a_i$ from a set of possible actions that has been pre-defined at the start. From the action, the scoring system will be updated, and the agent will receive a reward based on the change in to the scoring system.

A pseudocode of the DQN algorithm adopted from the paper is provided in Figure 2.

```
Initialize replay memory D to capacity N
Initialize action-value function Q with random weights θ
Initialize target action-value function Q̂ with weights θ⁻ = θ
For episode = 1, M do
    Initialize sequence s₁ = {x₁} and preprocessed sequence φ₁ = φ(s₁)
    For t = 1,T do
        With probability ε select a random action aₜ
        otherwise select aₜ = argmaxₐ Q(φ(sₜ),a; θ)
        Execute action aₜ in emulator and observe reward rₜ and image xₜ₊₁
        Set sₜ₊₁ = sₜ,aₜ,xₜ₊₁ and preprocess φₜ₊₁ = φ(sₜ₊₁)
        Store transition (φₜ,aₜ,rₜ,φₜ₊₁) in D
        Sample random minibatch of transitions (φⱼ,aⱼ,rⱼ,φⱼ₊₁) from D
```

$$y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$$

```
        Perform a gradient descent step on (yⱼ − Q(φⱼ,aⱼ; θ))²  with respect to the
        network parameters θ
        Every C steps reset Q̂ = Q
    End For
End For
```

**Figure 2** DQN Algorithm Pseudocode

*Exploration – To find out all the different possible states in the solution space.* Allows the agent to wander on its own to explore the space. After an action, compute the rewards and update the results into an action value function called *Q(s,a)* After 10,000 epochs of exploration, sample a mini-batch to train the Deep Neural Network Model to the update Q value.

*Exploitation – Utilize gained experience in order to obtain the optimal solution at the situation and attempts to maximize the rewards over time,* a.k.a. Greedy Approach. The downside of greedy approach is that it provides little or no exploratory potential. To solve this issue, we take a random action occasionally while choosing the optimal route.
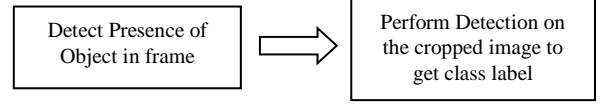
### D. Rewards

Rewards are almost exclusive to each case, but most applicable ones are:
  i. Distance to goal – More rewards for getting closer to the goal, penalize when the agent moves further away.
  ii. Object detected – Rewards are given whenever the agent finds something interesting or performs an interactive action that is relevant to the set objective.
  iii. Discovery of solution space – Some business use cases may require more exploration on new areas; in such cases the identification of new and potentially interesting areas can be rewarded.

### E. Computer Vision based Objectives

Most of the objectives of autonomous agents require the use of computer vision techniques. In cases such as monitoring using object detection, approaches such as two-step approach as suggested in [10] can be adopted. First step can be to identify the presence of the object in the frame, and crop the image using the bounding box coordinates. After enlarging the image,

a better prediction can be made to identify the exact object. This approach is suggested due to the possibility that the agent may not be in enough proximity to the object of interest, in such cases the two-step approach gives better results.



### F. Video Feed Transfer in Near Real Time

To improve the speed of transfer of the video frames from the agent to a server (for future analysis) individual image frames are compressed using ZStd algorithm from Facebook [11] and serialized and transferred across the network. The process is reversed for restoration of the original image. A typical jpg image of 900KB can be reduced to 20KB and transferred as shown in Figure 3.
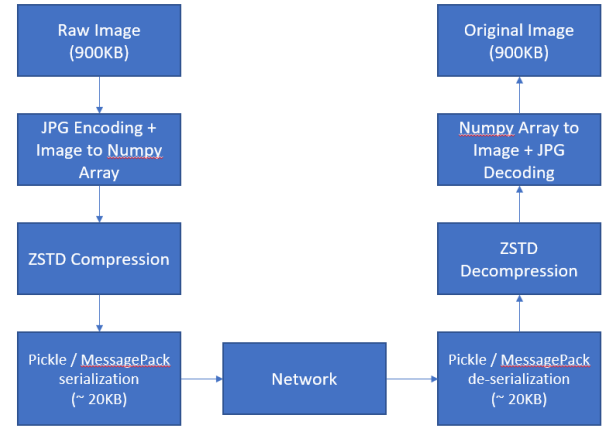


**Figure 3** Transferring image data from agent to server
The encoding-decoding process at both ends

### V. CASE STUDY – MONITORING WILDLIFE IN THE AFRICAN TERRAIN USING AN AUTONOMOUS UAV

As a case study, we focus on monitoring wildlife using a UAV. The simulation framework used is Microsoft's AirSim [6] which is an open source simulator for autonomous vehicles built on Unreal Engine. There are two options to select –i. quadcopters and ii. cars. African wildlife [9] was chosen as the environmental setting as the simulation focuses on wild animals. The objective is to detect as many animals as possible while navigating between two fixed points. [7][8]. Here we assume that the animals under observation do not migrate to far away distances (unlike birds) and therefore the environment is relatively static.

The following sub-sections describe the case specific solution approach, training and results of DQN for navigation and visualizing the route the drone takes during its flight.

Agent – The agent will be a drone (quadcopter) which has been pre-defined in the AirSim. It is a standard 4 rotor off the

shelf drone. In Unreal Engine, 4 cameras have been attached to it showcasing the different viewpoints of the camera. This is similar to an actual drone.

Functionalities of the drone includes
1) *4 Cameras (Left Center, Right Center, Center, Back and Bottom)*
2) *GPS Co-ordinates provided by Unreal Engine.*
3) *Velocity of drone provided by Unreal Engine*

Action – There will be 6 different actions available to the drone at each state. Since the drone moves in 3d space, we take into consider the three axes – x, y and z. The agent will choose one of the action available from this set which aims to maximize the Q-value. Each action will involve movement by a constant velocity factor.

1,2: Move in the positive, negative x-axis
3,4: Move in the positive, negative y-axis
5,6: Move in the positive, negative z-axis

Policy – The policy is influenced by the bottom camera, since the detection of animals results in more rewards.

Reward – The reward here is +100 whenever an animal is spotted by the bottom camera and detection by the object detection algorithm.

### A. Solution Approach



**Figure 4** Finding optimal path to maximize count of animals

The solution for this case study is to navigate the drone from point X (source) to Y (destination) and find the optimal route that covers maximum animals along the path.

### B. Training with Deep Q-Learning

The drone was trained on the environment for several episodes for days, the overall training was ~1.1M iterations. The graphs below provide the training metrics across iterations. Do note that the training has not converged and
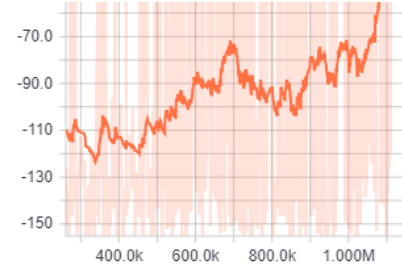


**Figure 5** Mean Action-Value (Q) vs Training Iterations

Initial action-value (Q) is very low due to random movements by the drone, a.k.a. the *exploration* phase. Eventually, the drone starts improving (negative value decreases) as it starts *exploiting* the familiar route.
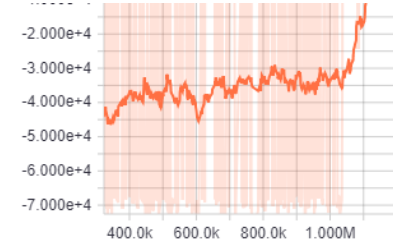


**Figure 6** Rewards per Training Iteration

The rewards are negative because of high penalty assigned for straying away from the path. As the iterations increase, the rewards improve as shown in Figure 5.
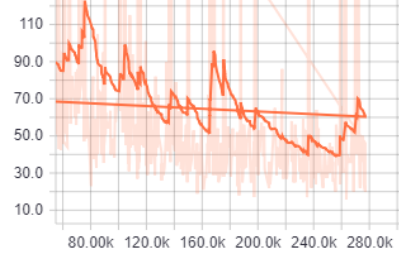


**Figure 7** Average Loss per minibatch

Figure 6 shows the average loss decreasing with every iteration, computed in terms of minibatch size of 4.

### C. Objective – Detection of wildlife using Computer Vision

Using the two-step approach similar to [10], the animals are first detected on the basis of an object captured is considered as present, the bounding boxes are cropped, and another detection algorithm identifies the species of animal from the enlarged image.
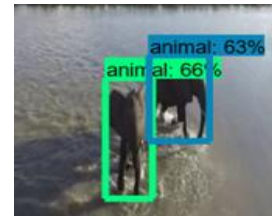


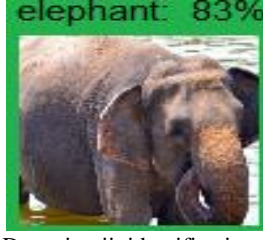**Figure 8** Object detection i. Presence of animal

**Figure 9** Object Detection ii. identification of specific animal

*D. Oversight of the route taken*

This sub-section describes the visualization of real-time route navigation from the drone as it traverses from start (star) to the destination (diamond).

After 1.1M iterations, the screenshot below shows the path traversed with an allowable distance deviation of 30m from the destination. It implies that the drone can land at any coordinate which is at most 30m away from the described destination.
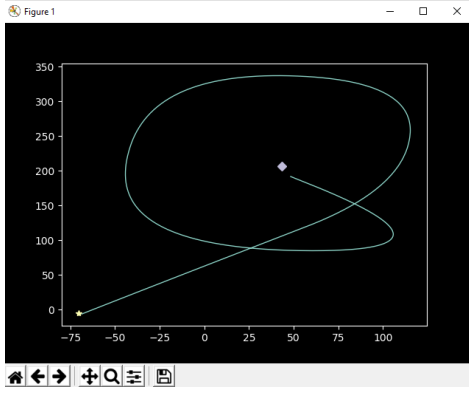


**Figure 10** Route visualization with acceptable distance deviation of 30m from destination

It is observed that when decreasing this deviation distance, the drone takes a longer time and covers wider route before landing. The screenshot below shows the route when the deviation distance is 5m.
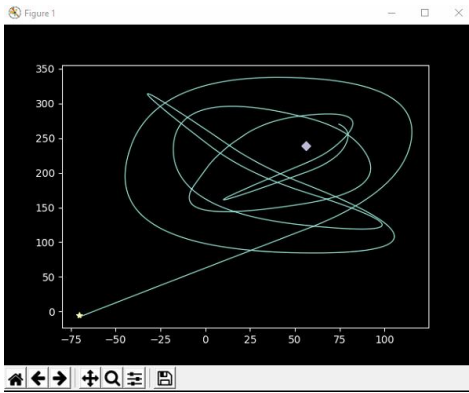


**Figure 11** Route visualization with acceptable distance deviation of 5m from destination

Therefore, based on the allowable deviation distance, we can fine tune the training time, making it converge sooner or later. The overall distance traveled also depends on the agent's

freedom of movement across dimensions – a drone has a lot more factors to consider due to its movement in the z-axis.
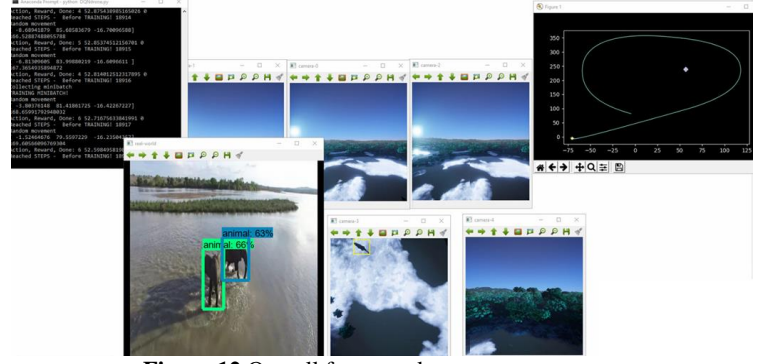


**Figure 12** Overall framework components

## VI. CHALLENGES AND FUTURE WORK

One of the existing challenges with reinforcement learning is the sheer amount of time it takes for training a decent model. Despite the environment space being bounded to restrict infinite possibilities, the large search space multiplied by the dimensions accessible by the agent makes the problem on an exponential scale. A future work can be to parallelly train multiple agents locally and aggregate the best features globally every few iterations (adopting a genetic algorithm approach) so that the solution converges faster.

Another work that can be attempted is to transfer learn the agent's model across environments and even attempt cross agent transfer and observe the results.

## VII. CONCLUSION

A generic conceptual framework for autonomous navigation was developed in an attempt to train an agent to move across in an environment and satisfy specified objectives by incentivizing with rewards based on performance. A case study was chosen – wildlife detection using drone in African terrain to showcase the implementation of the framework.

REFERENCES

[1] SUTTON, R. & BARTO, A. REINFORCEMENT LEARNING: AN INTRODUCTION (MIT PRESS, 1998)

[2] MNIH, V., KAVUKCUOGLU, K., SILVER, D., RUSU, A. A., VENESS, J., BELLEMARE, M. G & PETERSEN, S. (2015). HUMAN-LEVEL CONTROL THROUGH DEEP REINFORCEMENT LEARNING. *NATURE*, *518*(7540), 529.

[3] SILVER, D., HUANG, A., MADDISON, C. J., GUEZ, A., SIFRE, L., VAN DEN DRIESSCHE, G & DIELEMAN, S. (2016). MASTERING THE GAME OF GO WITH DEEP NEURAL NETWORKS AND TREE SEARCH. *NATURE*, *529*(7587), 484.

[4] VINYALS, O., EWALDS, T., BARTUNOV, S., GEORGIEV, P., VEZHNEVETS, A. S., YEO, M. & QUAN, J. (2017). STARCRAFT II: A NEW CHALLENGE FOR REINFORCEMENT LEARNING. *ARXIV PREPRINT ARXIV:1708.04782*.

[5] OPENAI RELEASE- LEARNING TO DEFEAT AMATEUR PLAYERS IN DOTA2 HTTPS://BLOG.OPENAI.COM/OPENAI-FIVE/

[6] AIRSIM: OPEN-SOURCE SIMULATOR FOR AUTONOMOUS VEHICLES BUILT ON UNREAL ENGINE HTTPS://GITHUB.COM/MICROSOFT/AIRSIM

[7] BONDI, E., FANG, F., HAMILTON, M., KAR, D., DMELLO, D., CHOI, J. & NEVATIA, R. (2018). SPOT POACHERS IN ACTION: AUGMENTING CONSERVATION DRONES WITH AUTOMATIC DETECTION IN NEAR REAL TIME. IAAI.

[8] BONDI, E., KAPOOR, A., DEY, D., PIAVIS, J., SHAH, S., HANNAFORD, R. & TAMBE, M. (2018). NEAR REAL-TIME DETECTION OF POACHERS FROM DRONES IN AIRSIM. IN *IJCAI* (PP. 5814-5816).

[9] AIRSIM-W: A SIMULATION ENVIRONMENT FOR WILDLIFE CONSERVATION WITH UAVS

[10] XULEI YANG ET AL (2018) 'DEEP LEARNING FOR PRACTICAL IMAGE RECOGNITION: CASE STUDY ON KAGGLE COMPETITIONS' IN KDD18

[11] ZSTANDARD REAL TIME COMPRESSION ALGORITHM HTTPS://FACEBOOK.GITHUB.IO/ZSTD/

[12] N. ZIMMERMAN ET AL., (2004) IMPLEMENTING A RULE-BASED SYSTEM TO REPRESENT DECISION CRITERIA FOR ON-ROAD AUTONOMOUS NAVIGATION, AAAI

[13] A. STEVENS ; M. STEVENS ; H. DURRANT-WHYTE (1995) "OXNAV": RELIABLE AUTONOMOUS NAVIGATION

[14] YIYOU SUN ; TONGHUA SU ; ZHIYING TU (2017) FASTER R-CNN BASED AUTONOMOUS NAVIGATION FOR VEHICLES IN WAREHOUSE

[15] CHAO WANG ; JIAN WANG ; XUDONG ZHANG ; XIAO ZHANG (2017) AUTONOMOUS NAVIGATION OF UAV IN LARGE-SCALE UNKNOWN COMPLEX ENVIRONMENT WITH DEEP REINFORCEMENT LEARNING