

Generalized Aircraft Collision Avoidance Through Deep Reinforcement Learning

ACAS1: Dharmesh Tarapore, Vincent Wahl, Kasim Patel, Shantanu Bobhate
dharmesh@bu.edu, vinwah@bu.edu, kasimp93@bu.edu, sbobhate@bu.edu

Abstract—The gradual integration of ADS-B¹ into the National Airspace System (NAS) has spurred research into its possible use in collision avoidance systems. While most new systems have shown to be much safer than the existing TCAS II (Traffic Alert and Collision Avoidance System), all of them make comparable assumptions about aircraft capabilities, thus restricting their applicability to highly specific classes of airplanes [?]. In this project, we develop a model derived from one such solution and generalize it to almost all powered aircraft by making conservative initial assumptions about their capabilities and then improving them by extrapolating from state action pairs. This modification allows us to provide a truly comprehensive collision avoidance system that can be used in most powered airplanes.

I. INTRODUCTION

The federally mandated Traffic Alert and Collision Avoidance System (TCAS II) for transport category aircraft has proven remarkably effective at averting mid-air collisions by providing pilots with timely alerts to resolve imminent threats. However, strong assumptions about aircraft capabilities made in the TCAS II logic prevent it from being used in general aviation, where the risk of a mid-air collision is significantly higher [?].

Newer collision avoidance systems handle the complexity induced by the steep variance in airplane dynamics by representing the problem as a Partially Observable Markov Decision Process (POMDP)[?]. Solving a POMDP for a given state-action pair (a, s) yields an action a' on the basis of a policy, π , and the corresponding reward, r , for choosing that action given the initial state. Each discrete state-action pair, its corresponding action, and reward is then stored in a policy table, which can be used to compare future input state-action pairs and output actions on the basis of a policy π that maximizes the expected utility for comparable state-action pairs found in the previously saved policy table. Selecting optimal actions thus requires storing millions of finely tuned policy tables, some of which have an average file size on the order of tens of gigabytes. This poses a serious problem for avionics onboard light aircraft, since budget constraints often restrict hardware memory and processing capabilities. Further, the granularity required to provide precise threat resolution policies grows exponentially as the POMDP learning progresses, eliminating the possibility of an online policy estimator that can refine its policy during actual usage [?].

¹ADS-B or automatic dependent surveillance-broadcast enables aircraft to determine its position using satellite navigation, which it then broadcasts periodically. ADS-B can serve as a supplement to secondary radar and is eventually slated to replace it.

Previous approaches have focused on evaluating the feasibility of using deep reinforcement learning to develop collision avoidance systems for Unmanned Aerial Systems (UAS) by approximating POMDP score tables [?]. Yet, all of the strategies explored so far have made strong assumptions about aircraft performance parameters, thus limiting their applicability to a particular class of aircraft. The approach documented by Julian, for example, makes strong assumptions about an aircraft's permissible turn rates to allow for abrupt maneuvers, which, while pragmatic on an unmanned aerial system, cannot extend to manned aircraft constrained by weight-induced structural limitations and delayed human response times [?]. Conversely, it also fails to exploit the enhanced performance capabilities of a manned aircraft, further restricting its utility. In this project, we extend the ideas discussed in [?] while making minimal assumptions about aircraft and sensor capabilities, opting instead to refine them for each individual aircraft on the basis of rational assumptions inferred from their state history. We do this to dynamically estimate a bound action space, \hat{A} , for an aircraft given an initial state-action pair describing its initial trajectory.

In particular, we begin with a limited set of assumptions about the available actions and modify that set on the basis of state information. This information is then input into a neural network to learn an approximation of the optimized score table. We evaluate its performance against TCAS II by comparing both algorithms' risk ratios² in a series of 15000 scripted encounters.

We further explore methods to improve performance through adaptive stress testing, a blackbox Monte Carlo Tree Search framework designed to maximize the likelihood of encounter trajectories that trigger near mid-air collisions.

II. RELATED WORK

Temizer et al. used POMDPs to develop a collision avoidance system for unmanned aircraft in [?]. Because POMDPs work with belief-states as opposed to exact states, it took the authors over 24 hours to generate an initial heuristic policy within acceptable bounds. Nonetheless, their policy generator outperformed TCAS II by 10%, thus demonstrating the potential superiority of POMDPs if solved efficiently. Advancements in deep reinforcement learning motivated Google's

²The risk ratio for a collision avoidance system, CAS , is given expressed as the probability of a near mid-air collision when using the collision avoidance system, $P(NMAC | CAS)$, divided by the probability of a near mid-air collision in the absence of any collision avoidance system, $P(NMAC)$. Smaller risk ratios across varied encounters usually signal better algorithms.

DeepMind group to combine deep reinforcement learning with heuristic search algorithms to develop an agent that defeated the European Go champion with a score of 5 - 0 [?]. Julian extended this idea to collision avoidance for unmanned aerial systems in [?], demonstrating policy approximations that were virtually indistinguishable from those produced by the uncompressed POMDP policy tables. We believe we can generalize this approach to develop collision avoidance strategies for most airplanes. As a global function approximator, an appropriately trained neural network is well poised to refine our initially conservative performance assumptions, thus relaxing the tight constraints imposed on our collision avoidance strategy's action space. This will allow us to output superior actions that resolve most near mid-air collision scenarios.

III. APPROACH

Similar to the approaches charted in [?] and [?], we expressed the problem as a POMDP, which comprises a state space S and an action space A . The robustness of applying POMDPs to uncertain state transitions and aircraft collision avoidance in particular, is described in detail in the book, *Decision Making Under Uncertainty* [?], which we will use to provide an overview of POMDPs. An agent in state $s \in S$ chooses an action $a \in A$ to receive a reward r and proceeds to state s' with probability $T(s'|a, s)$. The action taken is chosen on the basis of a policy π . An optimal policy maximizes the expected utility³. Given that we cannot have complete knowledge of the function describing the transition between states, we estimate the optimal values for a state action pair, $Q(s, a)$, by using samples of (s, a, r, s') and a learning parameter, α . This yields the modified Bellman equation:

$$Q(s, a) = Q(s, a) + \alpha * [R + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)]$$

We use this to model a scenario in which 2 aircraft, **ownship**, controlled by our policy π , and **intruder**—representing the threat, are flying along trajectories likely to end in a near mid-air collision. The Q value gives us the expected utility of taking action $a \in A$ while ownship is in state s , under the assumption that we will continue to pick the optimal policy in each iteration until the end of the encounter. The objective then becomes computing π such that given the states of both ownship and intruder, ownship will pick the optimal action, a . We alter this formulation slightly by encapsulating the action space in a set $\hat{A} \subseteq A$ such that for each aircraft equipped with our system, \hat{A} 's interval evolves (on the basis of state information) to ultimately span the widest possible subset of A . This helps us customize resolution advisories to fully exploit individual aircraft capabilities, drastically improving our model and the quality of its resolution advisories. In the rest of this section, we describe the model we used to simulate aircraft encounters, our neural network's inputs, outputs, loss function, and a summary of the code used to implement our solution. We conclude with by examining

³The expected utility is the sum of the immediate reward and the discounted utility of future states. Simplified, we value immediate rewards over future rewards and quantify this by assigning diminishing values to future rewards.

potential problems with our approach and future solutions to remedy them.

A. Encounter Modeling

The probability of a near mid-air collision (NMAC) under normal operating conditions is extremely low, thus making it difficult and time-consuming to learn optimal strategies by simulating normal aircraft trajectories. Further, simulated aircraft encounters need to have high fidelity and realistic encounter modeling to facilitate the best possible learning since unrealistic scenarios will do little to help the learning process. For instance, encounter scenarios in which aircraft are climbing while simultaneously banking are quite unrealistic. Further, we cannot assume that pilots will react instantly or correctly to any resolution advisories. Maneuvering counter to an advisory, even momentarily, can exacerbate the situation and demand aggressive corrective maneuvers to resolve. This makes simulating varied realistic aircraft encounters a monumental yet necessary task. Under guidance from the United States Department of Homeland Security and the Federal Aviation Administration, MIT's Lincoln Laboratory developed the *Collision Avoidance System Safety Assessment Tool* (CASSATT) framework [?]. CASSATT uses a 4 degree-of-freedom model to update aircraft state by applying the necessary airspeed acceleration, roll rate, and pitch rate to achieve realistic aircraft dynamics while assuming curvilinear motion with a zero-sideslip constraint⁴. CASSATT also includes an ICAO⁵ provided pilot response model to simulate pilot reactions to resolution advisories. We were able to procure Stanford University's *Reinforcement Learning Encounter Simulator* (RLES), an unclassified, open-source Julia implementation of CASSATT built at the Stanford Intelligent Systems Laboratory (SISL) [?]. Using RLES, we integrated our neural network, written in Python using the Keras and Tensorflow frameworks to evaluate its efficiency.

Work done by the ACAS2 team to improve aircraft encounter modeling using the *OpenSky* network further helped us relax the zero-sideslip constraint, thus allowing us to model transient dynamics with higher fidelity⁶ [?].

NB: It is worth mentioning here that arguably the strongest assumption made by this report is that **all encounters are correlated**. That is, we assume that of the 2 aircraft involved in an imminent NMAC, both aircraft have transponders that report their position, velocities, and accelerations, and that at least one of them is communicating with air traffic control. This is done because our simulator is written primarily with

⁴The zero-sideslip constraint ensures that no simulated aircraft turns are uncoordinated. This is a sensible assumption to make in most cases, but relaxing it allows for more aggressive maneuvers (within the aircraft's structural limitations), thus allowing us to resolve complicated encounter where there are no other alternatives

⁵The International Civil Aviation Organisation, or ICAO is a United Nations body that attempts to standardize aviation regulations around the world.

⁶Higher fidelity in a simulation offers us more control over the dynamic variables that describe an airplane's trajectory (described in more detail in the next section). This finer control lets us develop strategies that may contain maneuvers that comparable systems could not have considered.

the correlated encounter model developed by MIT's Lincoln Laboratory in mind, but also because the Federal Aviation Administration has legally mandated the use of ADS-B out transmitters on most aircraft seeking to fly in controlled airspace by 2020. We further assume, like all other collision avoidance systems, that our collision avoidance system is only invoked as a last line of defense: i.e. air traffic control has been unable to assist the pilot, ownship's pilots are unable to visually locate the intruder, and that the intruding aircraft is not communicating with either ATC or the pilot(s) of the ownship. Uncertainty stemming from uncooperative obstacles—aircraft with no transponders; birds; drones and other obstacles invisible to radar, will be investigated in the future.

B. Inputs

Our implementation (split between Julia and Python to facilitate unit conversions and develop a modular RLES interface)⁷, takes in the following inputs for 2 aircraft, AC1 and AC2:

- **Aircraft IDs** $ID1, ID2$: We designate $ID1$ and $ID2$ as the unique identifiers for the 2 aircraft.
- **Vertical Miss Distance** vmd : the difference in altitude between the two aircraft at the point of closest approach.
- **Horizontal Miss Distance** hmd : the horizontal range between the two aircraft at the point of closest approach.
- **Airspace class** A : This variable may take on one of four values: B, C, D, and O, indicating which class of airspace the encounter is in. The values B, C, and D correspond to the controlled airspace classes defined by the FAA. "O" represents "other airspace", which includes Class A, E, and G airspace in the United States.
- **Altitude Layer** L : Airspace is divided into five altitude layers. For their classification, we refer interested readers to [?]. The altitude layer for an encounter is determined by the altitude of AC1 at the **time of closest approach**, often abbreviated as TCA. The airspace class and altitude layer together can help us infer (at a minimum) ownship's capabilities. For instance, an aircraft flying in Class A airspace at 400 knots is likely capable of achieving a vertical velocity in the vicinity of about 5000 feet per minute, which is 10 times our initially assumed minimum vertical velocity rate of 500 feet per minute. Exigent scenarios caused by inadequate pilot action or other factors can be mitigated with greater certainty⁸ than in aircraft with relatively inferior capabilities.
- **Approach Angle** β : The heading of AC2 relative to AC1 at TCA (See Fig. 1).
- **Bearing** χ : The bearing of AC2 relative to AC1. Given β , hmd and χ , we can uniquely identify the lateral position and orientation of AC2 relative to AC1 at TCA [?].
- **Initial Airspeeds** v_1 and v_2 : Initial airspeeds of the two aircraft. We assume zero wind since aircraft close enough to be in an encounter situation are most likely within the same air mass.

⁷Most of these variables are described in greater detail in [?].

⁸by advising progressively more advanced and aggressive maneuvers

- **Accelerations** \dot{v}_1 and \dot{v}_2 : We assume constant acceleration given the short duration of the average encounter (~ 50 seconds).
- **Turn rate** $\dot{\psi}_1$ and $\dot{\psi}_2$: Turn rates may change every second.
- **Vertical accelerations** \dot{h}_1 and \dot{h}_2 : We assume these are constant vertical accelerations too.
- **Customized action space** A : This vector contains 4 actions that provide an estimate for the maximum vertical velocity and horizontal turn rates specific to ownship. These are estimated by the method explained above and are used to output actions whose magnitudes are less than or equal to the values outlined in this vector.

These 17 variables are provided to our Julia simulator interface (written in the file `ACASNNImpl.jl`), which condenses them into 5 vectors that are then input into our neural network implemented in the file `ACAS_NN.py`. The neural network then outputs [?].

We assume for the sake of simplicity that we are only resolving pairwise aircraft encounter scenarios, since analyzing multiple collision threats simultaneously was not possible given the time constraints. We also assume that all aircraft are capable of achieving a minimum vertical velocity of 500 feet per minute and turn rates of 20 degrees⁹ and improve and record these values for every aircraft whose state vectors demonstrate more advanced capabilities.

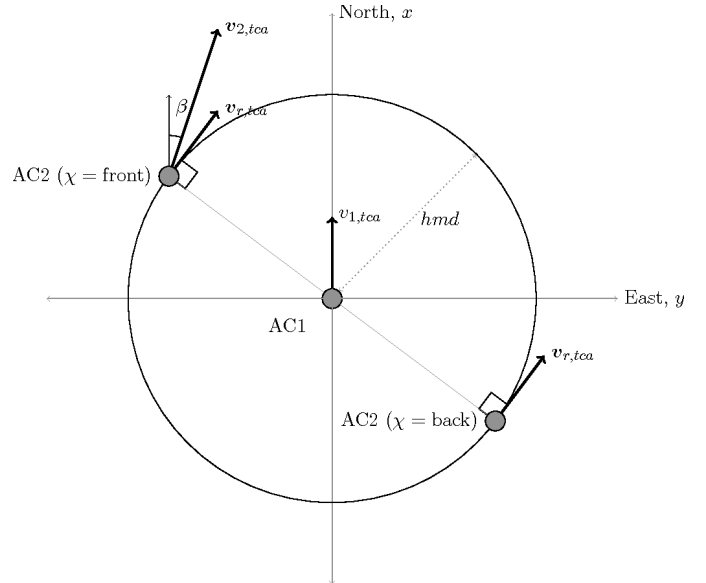


Fig. 1. Encounter geometry of 2 airplanes at their closest point of approach [?].

C. Outputs

The *TrainModel* method of our model uses state transition information to calculate a reward and stores it in the `acas_model.h5` file. Using this, the *Predict(state)* method

⁹These assumptions are consistent with current Air Traffic Controller expectations of aircraft operating in the United States National Airspace System (NAS).

outputs one of 5 actions: $\{a_1, a_2, a_3, a_4, a_5\}$ where a_1 and a_2 correspond to "Climb $h_i \leq h_{max}$ feet per minute" and "Descend $h_i \leq h_{max}$ feet per minute", respectively, while a_3 and a_4 correspond to "Turn left $\theta_i \leq \theta_{max}$ degrees" and "Turn right $\theta_i \leq \theta_{max}$ degrees", respectively. Here, h_{max} and θ_{max} represent the maximum possible vertical velocity and turn rate for the specific aircraft. These are calculated by drawing conclusions from state information when practicable. a_5 corresponds to the action *COC* or "Clear of Conflict". It describes a resolved encounter and allows ownship to move freely.

D. Reward

Our neural network showed suboptimal performance with large rewards and penalties, leading us to surmise that the discontinuities in our initial rewards (+100 for an action resulting in a resolved NMAC, -100 for an encounter that ended in a NMAC, -10 for unnecessary maneuvers) produced significant loss terms. A review of [?] encouraged us to smooth the reward function with a sigmoid function. Additionally, we adjusted our rewards scale estimate to make the absolute maximum reward be -1. Reward is thus calculated using the following metrics:

- $R_1 = -(1 + e^{(r_{sep}-r_{min})/C})^{-1}$ where hmd and vmd work together to describe r_{sep} , the distance of separation while r_{min} describes the minimum allowable separation (500 feet vertically and 200 feet horizontally¹⁰) and C is the smoothing applied to the step function.
- $R_2 = -0.0002\theta_i^2$ where θ_i represents the bank angle in degrees. This penalizes unnecessary maneuvering.
- $R_3 = -0.04$ if $\theta_i \neq COC \vee h_i \neq COC$, penalizing false positives.
- $R_4 = -0.0001h_i$, penalizing unnecessary vertical movement.

The final reward is given by $R = R_1 + R_2 + R_3 + R_4$.

E. Neural Network Implementation and Loss Function

Our initial model used a 3-hidden-layer neural network with rectified linear activations between the hidden layers. ReLu was chosen in favor of other activation functions given its relative computational efficiency. We chose 3 hidden layers to speed up computation and avoid overfitting¹¹. Instabilities inherent to reinforcement learning models using non-linear function approximators were well known and contained using the *experience replay* mechanism explained in [?]. We store the experiences $e_t = (s_t, a_t, r_t, s'_{t+1})$ at time-step t in memory. During the learning process, we apply Q-learning on batches of experiences sampled uniformly at random from the previously stored experiences to update the network parameters. A batch size of 128 was chosen with 84000 initial samples. Between batch update, we add 5 new samples while removing

the oldest 5 samples from memory. This is to handle the limited capacity of our replay memory. Removing the oldest actions ensures that newer actions are sampled more than older actions. Given the stochastic nature of state transitions in a POMDP, we use the expectation of the next state and reward to describe our loss function¹² as follows:

$$L(\theta) = \left[\mathbb{E}_{a' \in \mathcal{A}} (r + \gamma \max_{a'} Q(s', a', \theta^-)) - Q(s, a, \theta) \right]^2$$

It is worth noting that θ in the above function represents the neural network's parameters, not our turn rate input label.

F. Algorithm

The core algorithm is summarized below:

Algorithm 1 Generalized DQN ACAS Algorithm

```

1: procedure *(Input)
2:    $stringlen \leftarrow \text{length of } string$ 
3:    $i \leftarrow patlen$ 
4:   top:
5:     if  $i > stringlen$  then return false
6:      $j \leftarrow patlen$ 
7:     loop:
8:       if  $string(i) = path(j)$  then
9:          $j \leftarrow j - 1$ .
10:         $i \leftarrow i - 1$ .
11:        goto loop.
12:      close;
13:       $i \leftarrow i + \max(\delta_1(string(i)), \delta_2(j))$ .
14:      goto top.
```

TABLE I
SIMULATION PARAMETERS

Information message length	$k = 16000$ bit
Radio segment size	$b = 160$ bit
Rate of component codes	$R_{cc} = 1/3$
Polynomial of component encoders	$[1, 33/37, 25/37]_8$

IV. CONCLUSION

Gallia est omnis divisa in partes tres, quarum unam incolunt Belgae, aliam Aquitani, tertiam qui ipsorum lingua Celtae, nostra Galli appellantur. Gallos ab Aquitanis Garumna flumen, a Belgis Matrona et Sequana dividit. Horum omnium fortissimi sunt Belgae, propterea quod a cultu atque humanitate provinciae longissime absunt, minimeque ad eos mercatores saepe commeant atque ea quae ad effeminandos animos pertinent important, proximique sunt Germanis, qui trans Rhenum incolunt, quibuscum continenter bellum gerunt. Qua de causa Helvetii quoque reliquos Gallos virtute praecedunt, quod fere cotidianis proeliis cum Germanis contendunt, cum aut suis finibus eos prohibent aut ipsi in eorum finibus bellum gerunt. Eorum una, pars, quam Gallos obtinere dictum

¹⁰TCAS II thresholds often designate 500 feet of vertical and 100 feet horizontal separation as the minimum allowable separation. We chose 200 feet because it provides the pilot with more time to react, despite sacrificing potential resolutions to now untenable encounter scenarios.

¹¹Although we are now considering regularization as an alternative to reducing the number of layers

¹²Such an expectation can be approximated by simulating several trajectories and averaging the Q values and rewards provided by the next states.

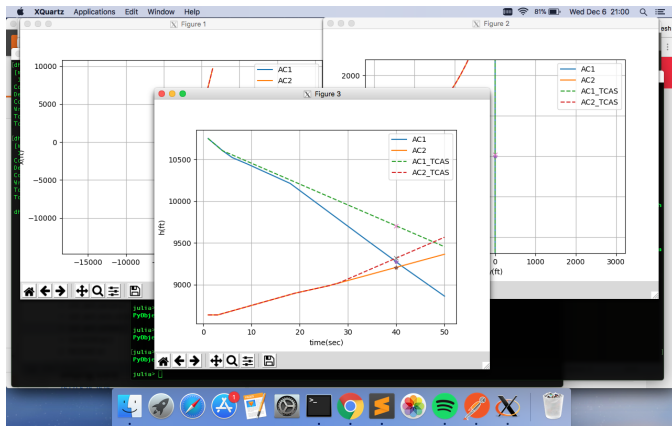


Fig. 2. A trajectory plot. For Vincent: We will add in a more visible plot later.

est, initium capit a flumine Rhodano, continetur Garumna flumine, Oceano, finibus Belgarum, attingit etiam ab Sequanis et Helvetiis flumen Rhenum, vergit ad septentriones. Belgae ab extremis Galliae finibus oriuntur, pertinent ad inferiorem partem fluminis Rheni, spectant in septentrionem et orientem solem.

V. CONCLUSION

This section summarizes the paper.

REFERENCES

- [1] J. Hagenauer, E. Offer, and L. Papke. Iterative decoding of binary block and convolutional codes. *IEEE Trans. Inform. Theory*, vol. 42, no. 2, pp. 429–445, Mar. 1996.
- [2] T. Mayer, H. Jenkac, and J. Hagenauer. Turbo base-station cooperation for intercell interference cancellation. *IEEE Int. Conf. Commun. (ICC)*, Istanbul, Turkey, pp. 356–361, June 2006.
- [3] J. G. Proakis. *Digital Communications*. McGraw-Hill Book Co., New York, USA, 3rd edition, 1995.
- [4] F. R. Kschischang. Giving a talk: Guidelines for the Preparation and Presentation of Technical Seminars. <http://www.comm.toronto.edu/frank/guide/guide.pdf>.
- [5] IEEE Transactions \LaTeX and Microsoft Word Style Files. <http://www.ieee.org/web/publications/authors/transjnl/index.html>