# Optimized and Trusted Collision Avoidance for Unmanned Aerial Vehicles using Approximate Dynamic Programming (Technical Report)

Zachary N. Sunberg, Mykel J. Kochenderfer, and Marco Pavone

*Abstract*—**Safely integrating unmanned aerial vehicles into civil airspace is contingent upon development of a trustworthy collision avoidance system. This paper proposes an approach whereby a parameterized resolution logic that is considered *trusted* for a given range of its parameters is adaptively tuned online. Specifically, to address the potential conservatism of the resolution logic with static parameters, we present a dynamic programming approach for adapting the parameters dynamically based on the encounter state. We compute the adaptation policy offline using a simulation-based approximate dynamic programming method that accommodates the high dimensionality of the problem. Numerical experiments show that this approach improves safety and operational performance compared to the baseline resolution logic, while retaining trustworthiness.**

## I. INTRODUCTION

As unmanned aerial vehicles (UAVs) move toward full autonomy, it is vital that they be capable of effectively responding to anomalous events, such as the intrusion of another aircraft into the vehicle's flight path. Minimizing collision risk for aircraft in general, and UAVs in particular, is challenging for a number of reasons. First, avoiding collision requires planning in a way that accounts for the large degree of uncertainty in the future paths of the aircraft. Second, the planning process must balance the competing goals of ensuring safety and avoiding disruption of normal operations. Many approaches have been proposed to address these challenges [1]–[9].

At present, there are two fundamentally different approaches to designing a conflict resolution system. The first approach focuses on inspiring confidence and trust in the system by making it as simple as possible for regulators and vehicle operators to understand by using hand-specified rules. Several algorithms that fit this paradigm have been proposed, and research toward formally verifying their safety-critical properties is underway [2]–[5]. In this paper, we will refer to such algorithms as *trusted resolution logics* (TRLs). TRLs typically have a number of parameters that determine how conservatively the system behaves, a feature that will be exploited later in this paper.

The second approach focuses on optimizing performance. This entails the offline or online computation of a "best" response action. Dynamic programming is widely used for this task [6]–[8]. Conflict resolution systems designed using this

Zachary Sunberg, Mykel Kochenderfer, and Marco Pavone are with the Department of Aeronautics & Astronautics, Stanford University, Stanford, CA 94305 {zsunberg, mykel, pavone}@stanford.edu

approach will be referred to as *directly optimized* systems. Unfortunately, even if a conflict resolution system performs well in simulation, government regulators and vehicle operators will often (and sometimes rightly) be wary of trusting its safety due to perceived complexity and unpredictability. Even in the best case, such a system would require expensive and time-consuming development of tools for validation as was the case for the recently developed replacement for the traffic alert and collision avoidance system (TCAS) [8].

In this paper, we propose a conflict resolution strategy that combines the strengths of these two approaches. The key idea is to use dynamic programming to find a policy that actively adjusts the parameters of a TRL to improve performance. If this TRL is trusted and certified for a range of parameter values, then the optimized version of the TRL should also be trusted and more easily certifiable.

We test the new approach in a scenario containing a UAV equipped with a perfect (noiseless) sensor to detect the state of an intruder and a simple TRL to resolve conflicts. This TRL, illustrated in Figure 1, determines a path that does not pass within a specified separation distance, $D$, of the intruder given that the intruder maintains its current heading (except in cases where no such path exists or when the TRL's heading resolution is too coarse to find such a path). To account for uncertainty in the intruder's flight path, if $D$ is fixed at a constant value, say $\bar{D}$, the value must be very large to ensure safety, but this may cause unnecessary departures from normal operation. To overcome this limitation, we compute an optimized policy $\tilde{\pi}$ that specifies a time varying separation distance, $D_t$, based on the encounter state. The goal is to ensure safety without being too conservative.

The problem of dynamically selecting $D_t$ can be formulated as a Markov decision process (MDP). Online solution of $\tilde{\pi}$ using an algorithm such as Monte Carlo Tree Search (MCTS) [10] would be conceptually straightforward, but would require significant computing power onboard the vehicle. In addition, it would be difficult and time-consuming to rigorously certify the implementation of MCTS due to its reliance on pseudo-random number generation. The key contribution of this paper is to devise an *offline* approach to compute $\tilde{\pi}$. Offline optimization of the policy is difficult because of the size of the state space of the MDP, which is the Cartesian product of the continuous state spaces of the UAV and the intruder. To overcome this challenge, we devise a value function approximation scheme that uses grid-based features that exploit the structure of the state space. This value function is optimized using
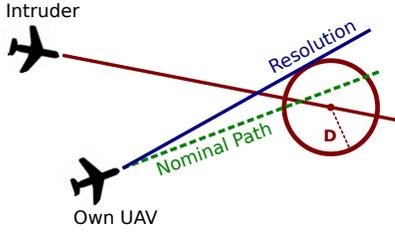
Fig. 1. Trusted Resolution Logic. The near mid air collision exclusion zone (circle with radius $D$) moves with the intruder through time, but is shown here only at the time of closest approach. The TRL (Algorithm 1) finds a straight trajectory close to the nominal path that avoids this zone.

simulation-based approximate dynamic programming (ADP). The policy is then encoded using an approximate post-decision state value function. With this post-decision value function, the UAV can easily extract the optimal action for the current state online by evaluating each possible action.

The remainder of the paper is organized as follows: Section II describes the MDP model of an encounter between two UAVs, Section III describes a solution approach based on approximate dynamic programming, and Section IV contains numerical results followed by conclusions in Section V.

## II. PROBLEM FORMULATION

The problem of avoiding an intruder using an optimized TRL is formulated as an MDP, referred to as the *encounter* MDP. An encounter involves two aerial vehicles flying in proximity. The first is the vehicle for which we are designing the resolution logic, which will be referred as the "own UAV" (or simply "UAV"). The second is the intruder vehicle, which may be manned or unmanned and will be referred to as the "intruder." Throughout the paper, the superscripts $^{(o)}$ and $^{(i)}$ refer to the own UAV and intruder quantities, respectively. Specifically, the encounter MDP is defined by the tuple $(S, A, T, R)$, which consists of

- The state space, $S$: The state of the encounter, $s$, consists of (1) the state of the own UAV, $s^{(o)}$, (2) the state of the intruder, $s^{(i)}$, and (3) a boolean variable $\mathrm{dev}$. The variable $\mathrm{dev}$ is set to true if the UAV has deviated from its nominal course and is included so that deviations at any point in time can be penalized equally. Collectively, the state is given by the triple

$$s = \left(s^{(o)}, s^{(i)}, \mathrm{dev}\right). \tag{1}$$

The state components $s^{(o)}$ and $s^{(i)}$ are specified in Section II-B. To model termination, $S$ also includes a termination state, denoted $s_{\mathrm{term}}$.
- The action space, $A$: The actions are the possible values for the separation distance $D \in \mathbb{R}_{\geq 0}$ used in the baseline TRL (see Figure 1).
- The state transition probability density function, $T$ : $S \times A \times S \rightarrow \mathbb{R}$: The value $T(s, D, s')$ is the probability density of transitioning to state $s'$ given that the separation parameter $D$ is used within the TRL at state $s$. This function is implicitly defined by a generative model that consists of a state transition function $F(\cdot)$ (described

in Section II-C) and a stochastic process $W$ (described in Section II-B).
- The reward, $R : S \times A \rightarrow \mathbb{R}$: The reward function, defined in Section II-D, rewards reaching a goal and penalizes near mid air collisions (NMACs), deviation from the nominal path, and time outside a goal region.

### A. Model assumptions

Two important simplifying assumptions were made for this initial study. First, the UAV and the intruder move only in the horizontal plane and at constant speed. Modeling horizontal maneuvers is necessary because UAVs will likely have to employ them in place of or in addition to the vertical maneuvers that current collision avoidance systems for manned aircraft such as TCAS rely on. This is due to both climb performance limitations and potential regulatory constraints such as the 500ft ceiling for small ($< 55$lb) UAVs in proposed Federal Aviation Administration rules [11]. Constraining altitude and speed simplifies exposition and reduces the size of the state space and hence the computational burden. Extensions to higher fidelity models (e.g., [12]) are possible and left for future research. A higher fidelity model would present a computational challenge, but perhaps not an insurmountable one. For example, the TRL could be extended to handle variable speed and altitude, but the policy that governs the TRL parameters could be optimized only on the most important dimensions of the model (e.g., the horizontal plane). See [13] for a similar successful example.

Second, the intruder dynamics are *independent* of the UAV's state; in other words, the intruder does not react to the flight path of the UAV. The cooperative setting where both the UAV and the intruder are equipped with a collision avoidance system (CAS) [6], [9], [14] is left for future research.

### B. Vehicle states and dynamics

This paper uses a very simple discrete-time model of an encounter between two aerial vehicles with time steps of duration $\Delta t$. Throughout this section, the superscript $^{(\cdot)}$ may be replaced by either $^{(o)}$ or $^{(i)}$. Both the UAV's and the intruder's states consist of the horizontal position $(x, y)$ and heading $\psi$, that is

$$s^{(o)} = \left(x^{(o)}, y^{(o)}, \psi^{(o)}\right), \quad s^{(i)} = \left(x^{(i)}, y^{(i)}, \psi^{(i)}\right). \tag{2}$$

The UAV and intruder also have similar dynamics. Both aircraft fly forward in the horizontal plane at constant speeds denoted by $v^{(\cdot)}$. They may turn at rates $\dot{\psi}^{(\cdot)}$ that remain constant over the simulation step. The following equations define the vehicle dynamics:

$$x_{t+1}^{(\cdot)} = \begin{cases} x_t^{(\cdot)} + v^{(\cdot)} \cos\left(\psi_t^{(\cdot)}\right) \Delta t & \text{if } \dot{\psi}^{(\cdot)} = 0 \\ x_t^{(\cdot)} + v^{(\cdot)} \frac{\sin\left(\psi_t^{(\cdot)} + w_t \Delta t\right) - \sin\left(\psi_t^{(\cdot)}\right)}{\dot{\psi}^{(\cdot)}} & \text{otherwise} \end{cases}$$

$$y_{t+1}^{(\cdot)} = \begin{cases} y_t^{(\cdot)} + v^{(\cdot)} \sin\left(\psi_t^{(\cdot)}\right) \Delta t & \text{if } \dot{\psi}^{(\cdot)} = 0 \\ y_t^{(\cdot)} - v^{(\cdot)} \frac{\cos\left(\psi_t^{(\cdot)} + \dot{\psi}^{(\cdot)} \Delta t\right) - \cos\left(\psi_t^{(\cdot)}\right)}{\dot{\psi}^{(\cdot)}} & \text{otherwise} \end{cases}$$

$$\psi_{t+1}^{(\cdot)} = \psi_t^{(\cdot)} + \dot{\psi}^{(\cdot)} \Delta t.$$

The intruder makes small random turns with

$$\dot{\psi}^{(\mathrm{i})} = w_t, \tag{3}$$

where $w_t$ is a stochastic disturbance. We let $W$ denote the stochastic process $\{w_t : t \in \mathbb{N}\}$. The random variables $w_t$ in $W$ are assumed independent and identically normally distributed with zero mean and a specified standard deviation, $\sigma_{\dot{\psi}}$. Subsequently, the intruder dynamics will be collectively referred to as $f^{(\mathrm{i})}\left(s_t^{(\mathrm{i})}, w_t\right)$.

The dynamics of the UAV are simplified conventional fixed wing aircraft dynamics with a single input, namely the roll angle $\phi^{(\mathrm{o})}$. We assume that the roll dynamics are fast compared to the other system dynamics, so that the roll angle $\phi^{(\mathrm{o})}$ may be directly and instantaneously commanded by the control system. This assumption avoids the inclusion of roll dynamics, which would increase the size of the state space.

$$\dot{\psi}_t^{(\mathrm{o})} = \frac{g \tan \phi_t^{(\mathrm{o})}}{v^{(\mathrm{o})}}, \tag{4}$$

where $g$ is acceleration due to gravity. The performance of the UAV is limited by a maximum bank angle

$$|\phi_t^{(\mathrm{o})}| \le \phi_{\max}. \tag{5}$$

The UAV dynamics will be collectively referred to as $f^{(\mathrm{o})}\left(s_t^{(\mathrm{o})}, \phi_t^{(\mathrm{o})}\right)$.

### C. UAV control system and transition function

The control system for the UAV consists of the TRL and a simple controller that commands a turn rate to track the course determined by the TRL. The TRL determines a heading angle, $\psi_{\mathrm{resolution}}^{(\mathrm{o})}$, that is (1) close to the heading to the goal, $\psi_{\mathrm{goal}}^{(\mathrm{o})}$, and (2) will avoid future conflicts with the intruder given that the intruder maintains its current heading as shown in Figure 1. In order to choose a suitable heading, many candidate headings, denoted $\psi_{\mathrm{cand}}^{(\mathrm{o})}$, are evaluated.

Given an initial state $s$, a candidate heading, $\psi_{\mathrm{cand}}^{(\mathrm{o})}$, for the UAV, and that both vehicles maintain their heading, the distance between the vehicles at the time of closest approach is a simple analytical function. Specifically, consider the distance $d\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}, \tau\right)$ between the vehicles $\tau$ time units in the future, i.e.,

$$d\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}, \tau\right) = \sqrt{\Delta x(\tau)^2 + \Delta y(\tau)^2}, \tag{6}$$

where

$$\Delta x(\tau) = x^{(\mathrm{i})} - x^{(\mathrm{o})} + \tau v^{(\mathrm{i})} \cos\left(\psi^{(\mathrm{i})}\right) - \tau v^{(\mathrm{o})} \cos\left(\psi_{\mathrm{cand}}^{(\mathrm{o})}\right),$$
$$\Delta y(\tau) = y^{(\mathrm{i})} - y^{(\mathrm{o})} + \tau v^{(\mathrm{i})} \sin\left(\psi^{(\mathrm{i})}\right) - \tau v^{(\mathrm{o})} \sin\left(\psi_{\mathrm{cand}}^{(\mathrm{o})}\right).$$

The minimum distance between the two vehicles is analytically found by setting the time derivative of $d\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}, \tau\right)$ to zero. Specifically, the time at which the vehicles are closest is given by

$$\tau_{\min}\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}\right) = \max\left\{\frac{a+b}{c-2d}, 0\right\}, \tag{7}$$

**Algorithm 1** Trusted Resolution Logic
___
**Input:** Encounter state $s$, desired separation distance $D$
**Output:** Resolution heading angle $\psi_{\mathrm{resolution}}^{(\mathrm{o})}$
  **function** TRL($s,D$)
    $\Psi \leftarrow \left\{\psi^{(\mathrm{o})} + n\pi/N : n \in \{-N, \ldots, N\}, N > 0\right\}$
    $D^* = \max_{\psi_{\mathrm{cand}}^{(\mathrm{o})} \in \Psi} d_{\min}\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}\right)$
    **if** $D^* < D$ **then**          $\triangleright$ conflict inescapable
      $\Psi \leftarrow \left\{\psi_{\mathrm{cand}}^{(\mathrm{o})} \in \Psi : d_{\min}\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}\right) = D^*\right\}$
      **return** $\underset{\psi_{\mathrm{cand}}^{(\mathrm{o})} \in \Psi}{\operatorname{argmin}} |\psi_{\mathrm{cand}}^{(\mathrm{o})} - \psi_{\mathrm{goal}}^{(\mathrm{o})}|$
    **else**
      $\Psi \leftarrow \left\{\psi_{\mathrm{cand}}^{(\mathrm{o})} \in \Psi : d_{\min}\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}\right) \ge D\right\}$
      **return** $\underset{\psi_{\mathrm{cand}}^{(\mathrm{o})} \in \Psi}{\operatorname{argmin}} |\psi_{\mathrm{cand}}^{(\mathrm{o})} - \psi_{\mathrm{goal}}^{(\mathrm{o})}|$
___

where

$$a := -v^{(\mathrm{i})} x^{(\mathrm{i})} \cos(\psi^{(\mathrm{i})}) - v^{(\mathrm{i})} y^{(\mathrm{i})} \sin(\psi^{(\mathrm{i})})$$
$$b := v^{(\mathrm{o})} x^{(\mathrm{i})} \cos(\psi_{\mathrm{cand}}^{(\mathrm{o})}) + v^{(\mathrm{o})} y^{(\mathrm{i})} \sin(\psi_{\mathrm{cand}}^{(\mathrm{o})})$$
$$c := v^{(\mathrm{o})\,2} + v^{(\mathrm{i})\,2} \cos^2(\psi^{(\mathrm{i})}) + v^{(\mathrm{i})\,2} \sin^2(\psi^{(\mathrm{i})})$$
$$d := v^{(\mathrm{o})} v^{(\mathrm{i})} (\cos(\psi^{(\mathrm{i})}) \cos(\psi_{\mathrm{cand}}^{(\mathrm{o})}) + \sin(\psi^{(\mathrm{i})}) \sin(\psi_{\mathrm{cand}}^{(\mathrm{o})})).$$

The minimum separation distance over all future time is then

$$d_{\min}\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}\right) := d\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}, \tau_{\min}\left(s, \psi_{\mathrm{cand}}^{(\mathrm{o})}\right)\right). \tag{8}$$

The TRL begins with a discrete set of potential heading values for the UAV. It then determines, for a desired separation distance $D$, which of those will not result in a collision given that the UAV and intruder maintain their headings. Finally, it selects the value from that set which is closest to $\psi_{\mathrm{goal}}^{(\mathrm{o})}$. The TRL is outlined in Algorithm 1.

Once the TRL has returned the desired heading, $\psi_{\mathrm{resolution}}^{(\mathrm{o})}$, a low-level controller determines the control input to the vehicle. We write this as

$$\dot{\psi}_t^{(\mathrm{o})} = c\left(s_t^{(\mathrm{o})}, \psi_{\mathrm{resolution}}^{(\mathrm{o})}\right), \tag{9}$$

where $c(\cdot)$ represents the low level controller.

The closed-looped dynamics for the state variables of the vehicles are then given by

$$s_{t+1}^{(\mathrm{i})} = f^{(\mathrm{i})}\left(s_t^{(\mathrm{i})}, w_t\right) \tag{10}$$
$$s_{t+1}^{(\mathrm{o})} = f^{(\mathrm{o})}\left(s_t^{(\mathrm{o})}, c\left(s_t^{(\mathrm{o})}, \mathrm{TRL}(s_t, D_t)\right)\right). \tag{11}$$

The UAV and the intruder dynamics are coupled *only* through the TRL.

The goal region that the UAV is trying to reach is denoted by $S_{\mathrm{goal}}$. This is the set of all states in $S$ for which $\left\|\left(x^{(\mathrm{o})}, y^{(\mathrm{o})}\right) - \left(x_{\mathrm{goal}}^{(\mathrm{o})}, y_{\mathrm{goal}}^{(\mathrm{o})}\right)\right\| \le D_{\mathrm{goal}}$, where $D_{\mathrm{goal}} > 0$ is a specified goal region radius, and $\left(x_{\mathrm{goal}}^{(\mathrm{o})}, y_{\mathrm{goal}}^{(\mathrm{o})}\right)$ is the goal center location. A near mid-air collision (NMAC) occurs at time $t$ if the UAV and intruder are within a minimum separation distance, $D_{\mathrm{NMAC}} > 0$, that is if $\left\|\left(x_t^{(\mathrm{o})}, y_t^{(\mathrm{o})}\right) - \left(x_t^{(\mathrm{i})}, y_t^{(\mathrm{i})}\right)\right\| \le D_{\mathrm{NMAC}}$. If the UAV reaches the goal region at some time $t$,

i.e., $s_t \in S_{\text{goal}}$, or if an NMAC occurs, the overall encounter state $s$ transitions to the terminal state $s_{\text{term}}$ and remains there. If the UAV performs a turn, dev is set to true because the vehicle has now deviated from the nominal straight path to the goal.

Let the state transition function, defined by (10), (11), and the special cases above, be denoted concisely as $F$ so that

$$s_{t+1} = F(s_t, D_t, w_t), \qquad (12)$$

where, as stated above, $D_t$ is the input to the system (the action in the MDP formulation).

### D. Reward

In this paper, we minimize two competing metrics. The first is the risk of an NMAC. As in previous studies (e.g., [8]), this aspect of performance is quantified using the *risk ratio*, the number of NMACs with the control system divided by the number of NMACs without the control system. The second metric is the probability of any deviation from the nominal path. This metric was chosen (as opposed to a metric that penalizes the magnitude of the deviation) because any deviation from the normal operating plan might have a large cost in the form of disrupting schedules, preventing a mission from being completed, or requiring manual human monitoring. The MDP reward function is designed to encourage a policy that performs well with respect to these goals.

Specifically, the total reward associated with an encounter is the sum of the stage-wise rewards throughout the entire encounter

$$\sum_{t=0}^{\infty} R(s_t, D_t). \qquad (13)$$

In order to meet both goals, the stage-wise reward is

$$R(s_t, D_t) := -c_{\text{step}} + r_{\text{goal}} \times \text{in\_goal}\left(s_t^{(\text{o})}\right)$$
$$- c_{\text{dev}} \times \text{deviates}(s_t, D_t)$$
$$- \lambda \times \text{is\_NMAC}(s_t), \qquad (14)$$

for positive constants $c_{\text{step}}$, $r_{\text{goal}}$, $c_{\text{dev}}$, and $\lambda$. The first term is a small constant cost accumulated in each step to push the policy to quickly reach the goal. The function in\_goal indicates that the UAV is within the goal region, so the second term is a reward for reaching the goal. The third term is a penalty for deviating from the nominal path. The function deviates returns 1 if the action will cause a deviation from the nominal course and 0 otherwise. It will only return 1 if the vehicle has not previously deviated and dev is false, so the penalty may only occur once during an encounter. Constants $c_{\text{step}}$, $r_{\text{goal}}$, and $c_{\text{dev}}$ represent relative weightings for the terms that incentivize a policy that reaches the goal quickly and minimizes the probability of deviation. Example values for these constants are given in Section IV. The fourth term is the cost for a collision. The weight $\lambda$ balances the two performance goals. We heuristically expect there to be a value of $\lambda$ for which the solution to the MDP meets the desired risk ratio if it is attainable. Bisection, or even a simple sweep of values can be used to find a suitable value, and this method has been used previously to analyze the performance of aircraft collision avoidance systems [6], [7].

### E. Problem statement

The problem we consider is to find a feedback control policy $\pi^* : S \to A$, mapping an encounter state $s_t$ into a separation distance $D_t$, that maximizes the expected reward (13) subject to the system dynamics (12):

$$\begin{aligned} \underset{\pi}{\text{maximize}} \quad & E\left[\sum_{t=0}^{\infty} R(s_t, \pi(s_t))\right] \\ \text{subject to} \quad & s_{t+1} = F(s_t, \pi(s_t), w_t), \end{aligned} \qquad (15)$$

for all initial states $s_0 \in S$. In the next section, we present a practical approach to problem (15) that uses approximate dynamic programming to find a suboptimal policy $\tilde{\pi}$.

## III. SOLUTION APPROACH

The solution approach for problem (15) is an approximate dynamic programming algorithm called approximate value iteration [16]. The value function, $V$, represents the expected value of the future reward given that the encounter is in state $s$ and an optimal policy will be executed in the future. We approximate $V$ with a linear architecture of the form

$$\tilde{V}(s) = \beta(s)^{\top}\theta, \qquad (16)$$

where the feature function $\beta$ returns a vector of $N_{\beta}$ feature values, and $\theta \in \mathbb{R}^{N_{\beta}}$ is a vector of weights [16]. At each step of value iteration, the weight vector $\theta$ is fitted to the results of a large number of single-step simulations by solving a linear least-squares problem. After value iteration has converged, $\tilde{V}$ is used to compute a linear approximation of the *post-decision state* value function, $\tilde{V}_q$. The policy is extracted online in real time by selecting the action that results in the post-decision state that has the highest value according to $\tilde{V}_q$. The choice of working with post-decision states will be discussed in Section III-B.

### A. Approximate value iteration

The bulk of the computation is carried out offline before vehicle deployment using simulation. Specifically, the first step is to estimate the optimal value function for problem (15) using value iteration [16]. On a continuous state space, the Bellman operator used in value iteration cannot be applied for each of the uncountably infinite number of states, so an approximation must be used. In this paper we adopt projected value iteration [16], which uses a finite number of parameters to approximate the value function. Each successive approximation, $\tilde{V}_k$, is the result of the Bellman operation projected onto a linear subspace, that is

$$\tilde{V}_{k+1}(s) = \Pi\mathcal{B}[\tilde{V}_k](s), \qquad (17)$$

where $\mathcal{B}$ is the Bellman operator, and $\Pi$ is a Euclidean projection onto the linear subspace $\Phi$ spanned by the $N_{\beta}$ basis functions (see [16] for a detailed discussion of this approach).

To perform the approximate value iteration (17), we resort to Monte Carlo simulations. Specifically, for each iteration, $N_{\text{state}}$ states are uniformly randomly selected. If the states lie within the grids used in the feature function (see Section III-D) the sample is "snapped" to the nearest grid point to prevent

approximation errors due to the Gibbs phenomenon [17]. At each sampled state $s^{[n]}$, $n = 1, \ldots, N_{\text{state}}$, the stage-wise reward and the expectation of the value function are evaluated for each action $a$ within a discrete approximation of the action space $A$, denoted by $\tilde{A}$. The expectation embedded in the Bellman operator is approximated using $N_{\text{EV}}$ *single-step* intruder simulations, each with a randomly generated noise value, $w_m$, $m = 1, \ldots, N_{\text{EV}}$. However, since the UAV dynamics are deterministic, only one UAV simulation is needed. The maximum over $\tilde{A}$ is stored as the $n$th entry of a vector $v_{k+1}$:

$$v_{k+1}[n] := $$
$$\max_{D \in \tilde{A}} \left\{ R(s^{[n]}, D) + \frac{1}{N_{\text{EV}}} \sum_{m=1}^{N_{\text{EV}}} \beta(F(s^{[n]}, D, w_m))^\top \theta_k \right\},$$

for $n = 1, \ldots, N_{\text{state}}$. Here $v_{k+1}$ provides an approximation to the (unprojected) value function. To project $v_{k+1}$ onto $\Phi$, we compute the weight vector $\theta_{k+1}$ by solving the least-squares optimization problem

$$\theta_{k+1} = \operatorname*{argmin}_{\theta \in \mathbb{R}^{N_\beta}} \sum_{n=1}^{N_{\text{state}}} \left( \beta \left( s^{[n]} \right)^\top \theta - v_{k+1}[n] \right)^2. \tag{18}$$

Iteration is terminated after a fixed number of steps, $N_{VI}$, and the resulting weight vector, denoted $\theta$, is stored for the next processing step (Section III-B).

### B. Post decision value function extraction

For reasons discussed in Section III-C, our second step is to approximate a value function, $V_q$, defined over post-decision states [16], [18]. A post-decision state, $q$, is a state in $S$ made up of the own UAV state and dev *at one time step into the future* and the intruder state *at the current time*, that is

$$q_t := \left( s_{t+1}^{(\text{o})}, s_t^{(\text{i})}, \text{dev}_{t+1} \right). \tag{19}$$

Correspondingly, let $g : S \times A \to S$ be the function that maps the current state and action to the post-decision state. In other words, function $g(s_t, D_t)$ returns $q_t$ consisting of

$$s_{t+1}^{(\text{o})} = f^{(\text{o})} \left( s_t^{(\text{o})}, c \left( s_t^{(\text{o})}, TRL(s_t, D_t) \right) \right)$$
$$s_t^{(\text{i})} = s_t^{(\text{i})}$$
$$\text{dev}_{t+1} = \max\{\text{dev}, \text{deviates}(s_t, D_t)\}. \tag{20}$$

The post-decision value function approximation, $\tilde{V}_q$, is computed as follows: Let $h : S \times \mathbb{R} \to S$ be a function that returns the next encounter state given the post decision state and intruder heading noise value, that is $h(q_t, w)$ returns $s_{t+1}$ consisting of

$$s_{t+1}^{(\text{o})} = s_{t+1}^{(\text{o})} \tag{21}$$
$$s_{t+1}^{(\text{i})} = f^{(\text{i})}(s_t^{(\text{i})}, w) \tag{22}$$
$$\text{dev}_{t+1} = \text{dev}_{t+1}, \tag{23}$$

where $\left( s_{t+1}^{(\text{o})}, s_t^{(\text{i})}, \text{dev}_{t+1} \right)$ are the components of $q_t$.

The post-decision value function, $V_q$, is, in terms of the value function, $V$,

$$V_q(q) = E\left[V\left(h\left(q, w\right)\right)\right], \tag{24}$$

where $w$ denotes, as usual, a random variable with Gaussian normal density. The approximation, $\tilde{V}_q$, is of the linear form

$$\tilde{V}_q(q) = \beta(q)^\top \theta^q, \tag{25}$$

where $\beta(q)$ is the feature vector for post-decision states $q \in S$, and $\theta^q \in \mathbb{R}^{N_\beta}$ is the corresponding weight vector.

To calculate the weight vector, $N_q$ post decision states are randomly selected using the same method as described in Section III-A and are denoted $q^{[n]}$, $n = 1, \ldots, N_q$. For each sample $q^{[n]}$, the expectation in (24) is approximated using $N_{\text{EV}}$ single-step simulations. The results are used to solve a least squares optimization problem

$$\theta^q = \operatorname*{argmin}_{\theta \in \mathbb{R}^{N_\beta}} \sum_{n=1}^{N_q} \left( \beta \left( q^{[n]} \right)^\top \theta - v^q[n] \right)^2, \tag{26}$$

where

$$v^q[n] := \frac{1}{N_{\text{EV}}} \sum_{m=1}^{N_{\text{EV}}} \beta \left( h \left( q^{[n]}, w_m \right) \right)^\top \theta, \tag{27}$$

where $w_m$, $m = 1, \ldots, N_{\text{EV}}$, is sampled from a random variable in $W$.

### C. Online policy evaluation

The first two steps (explained, respectively, in Sections III-A and III-B) are performed offline. The last step, namely policy evaluation, is performed online. Specifically a suboptimal control at any state $s$ is computed as

$$\tilde{\pi}(s) = \operatorname*{argmax}_{D \in \tilde{A}} \tilde{V}_q \left( g(s, D) \right). \tag{28}$$

Since $g$ (defined in (20)) is a *deterministic* function of a state-action pair, this calculation does not contain any computationally costly or difficult-to-certify operations such as estimating an expectation.

Two comments regarding the motivation for using the post-decision value function are in order. First, if the value functions could be exactly calculated, the post-decision state approach would be equivalent to the more common state-action value function approach, wherein a control is computed by solving

$$\pi(s) = \operatorname*{argmax}_{D \in \tilde{A}} Q(s, D). \tag{29}$$

Here, $Q(s, a)$ is the state-action value function representing expected value of taking action $a$ in state $s$ and then following an optimal policy [16], [18]. One can readily show that, in the exact case, $Q(s, D) = V_q \left( g(s, D) \right)$. However, for this problem, post-decision states are much less susceptible to approximation errors. This is primarily due to the fact that it is difficult to specify suitable features in the state-action domain. Since $g$ is highly nonlinear (indeed it is discontinuous because of the TRL), a grid interpolation approximation scheme (see Section III-D) is not suitable for approximating $Q$. However, by evaluating $g$ online at the current state in Equation (28),
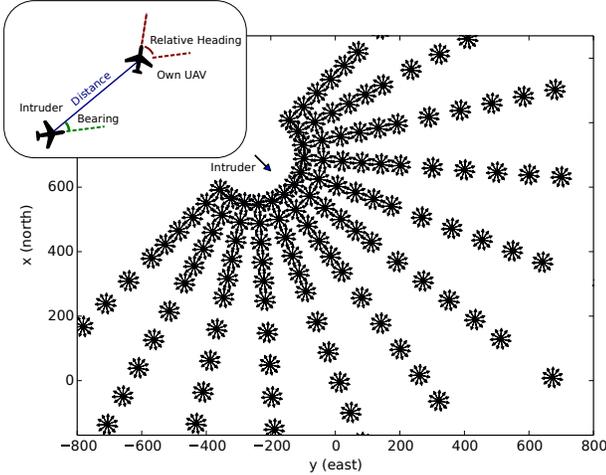
Fig. 2. Intruder interpolation grid for $\beta_{\text{intruder}}$, visualized with the intruder at $(700\,\text{m},\ -250\,\text{m})$ at heading 135. The top left inset shows the variables used. In the main plot, each of the small arrows represents a grid point.
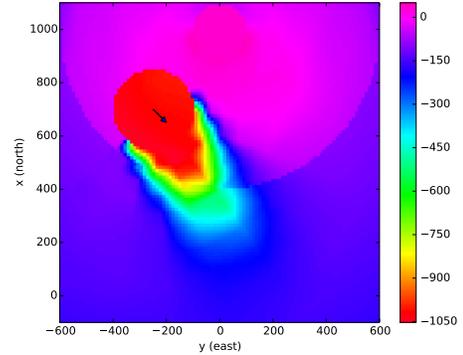


Fig. 3. Two-dimensional slice of the approximate optimal value function. Each pixel on this image represents the value function evaluated with the intruder at that position pointed directly north. The intruder is located at $(700\,\text{m}, -250\,\text{m})$ pointed at heading 135 as indicated by the arrow. The goal is at $(1000\,\text{m}, 0)$.

the difficulties with nonlinearity are avoided since $V_q$ is well-approximated by grid interpolation.

Second, the post decision value function is *not* used for the value iteration portion of the offline solution as the expectation estimate in Equation (27) would require $N_{EV}$ simulations of both the own UAV and the intruder dynamics and, therefore, would be more computationally demanding.

### D. Selection of features

The primary value function approximation features are the interpolation weights for points in a grid [19]. A grid-based approximation is potentially inefficient compared to a small number of global features (e.g. heading, distance, and trigonometric functions of those variables). However, it is well known that the function approximation used in value iteration must have suitable convergence properties [20] in addition to approximating the final value function. Indeed, we experimented with a small number of global features, but were unable to achieve convergence and resorted to using a grid. Since a grid defined over the entire six dimensional encounter state space with a reasonable resolution would require far too many points to be computationally feasible, the grid must be focused on important parts of the state space.

Our strategy is to separate features into two groups (along with a constant), specifically

$$\beta(s) = [\beta_{\text{intruder}}(s^{(\text{o})} - s^{(\text{i})}), \beta_{\text{goal}}(s^{(\text{o})}), 1]. \quad (30)$$

The first group, $\beta_{\text{intruder}}$, captures the features corresponding to a near midair collision and is a function of only the position and orientation of the UAV relative to the intruder. The second group, $\beta_{\text{goal}}$, captures the value of being near the goal and is a function of only the UAV state.

Since the domain of $\beta_{\text{intruder}}$ is only three dimensional and the domain of $\beta_{\text{goal}}$ is only two dimensional, relatively fine interpolation grids can be used for value function approximation without requiring a prohibitively large number of features. The $\beta_{\text{intruder}}$ feature group consists of a NMAC indicator function

and interpolation weights for a grid (Figure 2) with nodes at regularly spaced points along the following three variables: (1) the distance between the UAV and the intruder, (2) the bearing from the intruder to the UAV, and (3) the relative heading between the vehicles. The $\beta_{\text{goal}}$ vector consists of a goal indicator function, the distance between the UAV and the goal, and interpolation weights for a grid with nodes regularly spaced along the distance between the UAV and the goal and the absolute value of the bearing to the goal from the UAV. The total number of features is $N_\beta = 1813$.

Figure 3 shows a single two-dimensional "slice" of the full six-dimensional optimized value function. As expected, there is a low value region in front of and to the south of the intruder and an increase in the value near the goal.

Figure 4 shows a slice of an optimized policy. When multiple actions result in the same post-decision state value, the least conservative action is chosen, so the policy yields the lowest value of $D$ (500ft $\approx$ 152.4m) on most of the state space. Because the own UAV is pointed north, the policy is conservative in a region in front of and to the south of the intruder. The band corresponding to small $D$ that stretches across the middle of the conservative region (from $(700\,\text{m}, -250\,\text{m})$ to $(-100\,\text{m}, 200\,\text{m})$) is present because all values of $D$ result in the same post decision value.

## IV. RESULTS

This section presents results from numerical experiments that illustrate the effectiveness of our approach. The experiments are designed to compare the three approaches discussed in Section I: the "static TRL" approach, the "directly optimized" approach, and the newly proposed "optimized TRL" approach. Each is represented by a different control law that specifies a turn rate to the own UAV. Specifically, the static TRL law uses the TRL described in Algorithm 1 with a constant value for the separation distance (denoted $\bar{D}$). To represent the directly optimized approach, an optimized value function approximation and policy are generated using the method described in Section III, except that the policy *directly determines the turn rate*, $\dot{\psi}_t^{(\text{o})}$, *instead of specifying $D$*. The action
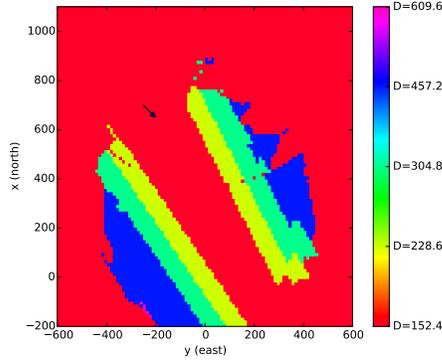
Fig. 4. Two-dimensional slice of the optimized policy for TRL parameter $D$. Each pixel in the image is the policy evaluated with the own UAV at that location pointed directly north. The intruder is located at $(700\,\text{m}, -250\,\text{m})$ and pointed at heading 135 as indicated by the arrow. Each color represents a different value for $D$ indicated by the number in the key on the right.

TABLE I
PARAMETERS FOR NUMERICAL EXPERIMENTS

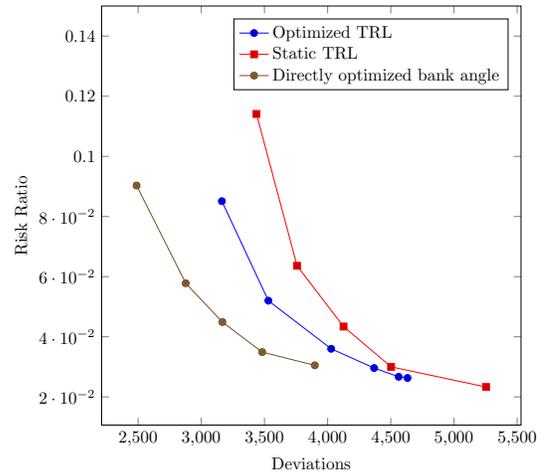| Description | Symbol | Value |
|---|---|---|
| Own UAV speed | $v^{(o)}$ | 30 m/s |
| Maximum own UAV turn rate | $\dot{\psi}_{\max}$ | 18.7°/s |
| Intruder speed | $v^{(i)}$ | 60 m/s |
| Intruder turn rate standard deviation | $\sigma_{\dot{\psi}}$ | 10°/s |
| Near mid air collision radius | $D_{\text{NMAC}}$ | 500 ft |
| Step cost | $c_{\text{step}}$ | 1 |
| Reward for reaching goal | $r_{\text{goal}}$ | 100 |
| Cost for deviation | $c_{\text{dev}}$ | 100 |
| Step simulations for expectation estimate | $N_{EV}$ | 20 |
| Single step simulations per round of value iteration (optimized TRL) | $N_{\text{state}}$ | 10,000 |
| Single step simulations per round of value iteration (directly optimized) | $N_{\text{state}}$ | 50,000 |
| Number of value iteration rounds | $N_{VI}$ | 35 |
| Single step simulations for post decision value function extraction | $N_q$ | 50,000 |

space for this policy is $\{-\dot{\psi}_{\max}, -\dot{\psi}_{\max}/2, 0, \dot{\psi}_{\max}/2, \dot{\psi}_{\max}\}$. The third optimized TRL law works as described in Sections II and III, dynamically assigning a value of $D$ to the TRL based on the state and using $c(\cdot)$ to assign the turn rate. The action space for the underlying $D$ policy is $\tilde{A} = \{D_{\text{NMAC}}, 1.5D_{\text{NMAC}}, 2D_{\text{NMAC}}, 3D_{\text{NMAC}}, 4D_{\text{NMAC}}\}$.

The parameters used in the numerical trials are listed in Table I. The control policies are evaluated by executing them in a large number of complete (from $t = 0$ to the end state) encounter simulations. The same random numbers used to generate intruder noise were reused across all collision avoidance strategies to ensure fairness of comparisons. In each of the simulations, the own UAV starts pointed north at position $(0, 0)$ in a north-east coordinate system with the goal at $(1000\text{m}, 0)$.

The evaluation simulations use the same intruder random turn rate model with standard deviation $\sigma_{\dot{\psi}}$ that was used for value iteration. A robustness study using different models is not presented here, but previous research [13] suggests that this method will offer good performance when evaluated against both a range of noise parameters and structurally different models. The intruder initial position is randomly generated between $800\,\text{m}$ and $1500\,\text{m}$ from the center point of the encounter area at $(500\text{m}, 500\text{m})$ with an initial heading that is within $135°$ of the direction from the initial position to the center point.

The conservativeness of each control law is characterized by counting the number of deviations from the nominal path in 10,000 simulations. Of these simulations, 1009 result in a NMAC if the UAV follows its nominal path, but for most a deviation would not be necessary to avoid the intruder.

The risk ratio is estimated using a separate set of 10,000 simulations. Each of these simulations has an initial condition in the same region described above, but initial conditions and noise trajectories are chosen by filtering random trials so that each of the simulations *will result in a NMAC if the own UAV follows its nominal path*. The risk ratio estimate for a policy is simply the fraction of these simulations that result in a NMAC when the policy is executed on them.



Fig. 5. Control policy comparison. The horizontal axis variable is the number of deviations in 10,000 encounter simulations. The values of $\lambda$ used to generate the datapoints are 100, 316, 1000, 3160, $10^4$, and $3.16 \times 10^4$ for the optimized TRL policy and 300, 500, 700, 1000, and 1500 for the directly optimized approach. The values of $\bar{D}$ for the static TRL policy are $250\,\text{m}$, $300\,\text{m}$, $350\,\text{m}$, $400\,\text{m}$, and $500\,\text{m}$.

Figure 5 shows the Pareto optimal frontiers for the different control laws. Each curve is generated by using various values of $\lambda$ in the reward function of the MDP, or by using various static values of $\bar{D}$ in the static TRL case. It is clear that optimization improves the performance of the TRL. For example, if the desired risk ratio is $5\%$, interpolation between data points suggests that the optimized TRL policy will cause approximately $10\%$ fewer deviations than the static TRL policy.

The directly optimized approach performs better than both TRL policies. This is not unexpected since the directly optimized policy is not limited to trajectories that the TRL deems to be safe. However, precisely because the directly optimized policy is not limited to trajectories guaranteed safe by hand-specified rules, it is much more difficult to convince people to trust it. These results thus allow us to estimate the performance *price* of using a trusted resolution logic rather than an untrusted one. Since the optimized TRL causes approximately

18% more deviations than the directly optimized policy at a risk ratio of 0.05, we may say the price of the trustworthy properties of the TRL at a desired risk ratio of 0.05 is an 18% increase in deviations. The new optimized TRL approach has reduced this from a price of a 31% increase in deviations without optimization.

The Julia code used for these experiments is available at `https://github.com/StanfordASL/UASEncounter`.

## V. Conclusion

This paper presented a method for improving the performance of trusted conflict resolution logic for an unmanned aerial vehicle through approximate dynamic programming. A linear value function approximation based on interpolation grids and other features produces policies that improve the performance of the resolution logic without undermining the trust placed in it. Specifically, simulation experiments show that this optimization approach is able to decrease the number of deviations from the nominal path without increasing the risk ratio. Comparison with a directly optimized approach that is more difficult to trust can help to quantify the price of the gain in trustworthiness that comes from using the trusted resolution logic. The optimization approach proposed here is capable of reducing that price.

There are several directions for future work on this problem. The current formulation models uncertainty in the intruder's actions, but sensor uncertainty should also be taken into account. If sensor uncertainty is added to the problem, it becomes a partially observable Markov decision process. One promising solution approach for the partially observable version of this problem is Monte Carlo value iteration [7]. This approach has been applied to UAV collision avoidance before [7], but it would be difficult to certify as a direct control system. Monte Carlo value iteration used in conjunction with a TRL has potential to handle both uncertainty in intruder behavior and sensor measurement noise while being easily certifiable. Also, as mentioned in Section II-A, there is potential to expand this work to more complicated scenarios with multiple intruders, intruders equipped with a CAS, higher fidelity models, or different sets of assumptions.

## References

[1] James K. Kuchar and Lee C. Yang, "A review of conflict detection and resolution modeling methods," *IEEE Transactions on Intelligent Transportation Systems*, vol. 1, no. 4, pp. 179–189, 2000.
[2] R. Bach, C. Farrell, and H. Erzberger, "An algorithm for level-aircraft conflict resolution," NASA, Tech. Rep. CR-2009-214573, 2009.
[3] G. Hagen, R. Butler, and J. Maddalon, "Stratway: A modular approach to strategic conflict resolution," in *AIAA Aviation Technology, Integration, and Operations (ATIO) Conference*, 2011.
[4] H. Herencia Zapana, J.-B. Jeannin, and C. Muñoz, "Formal verification of safety buffers for sate-based conflict detection and resolution," in *International Congress of the Aeronautical Sciences*, 2010.
[5] A. Narkawicz, C. Muñoz, and G. Dowek, "Provably correct conflict prevention bands algorithms," *Science of Computer Programming*, vol. 77, no. 10–11, pp. 1039–1057, 2012.
[6] M. J. Kochenderfer and J. P. Chryssanthacopoulos, "Robust airborne collision avoidance through dynamic programming," MIT Lincoln Laboratory, Tech. Rep. ATC-371, 2011.
[7] H. Bai, D. Hsu, M. J. Kochenderfer, and W. S. Lee, "Unmanned aircraft collision avoidance using continuous-state POMDPs," in *Robotics: Science and Systems*, 2012.
[8] J. E. Holland, M. J. Kochenderfer, and W. A. Olson, "Optimizing the next generation collision avoidance system for safe, suitable, and acceptable operational performance," *Air Traffic Control Quarterly*, vol. 21, no. 3, pp. 275–297, 2013.
[9] E. J. Rodríguez Seda, "Decentralized trajectory tracking with collision avoidance control for teams of unmanned vehicles with constant speed," in *American Control Conference*, 2014.
[10] A. Couëtoux, J.-B. Hoock, N. Sokolovska, O. Teytaud, and N. Bonnard, "Continuous upper confidence trees," in *Learning and Intelligent Optimization*. Springer, 2011, pp. 433–445.
[11] Federal Aviation Administration, "Small UAS notice of proposed rulemaking," February 2015. [Online]. Available: https://www.faa.gov/regulations_policies/rulemaking/recently_published/media/2120-AJ60_NPRM_2-15-2015_joint_signature.pdf
[12] M. J. Kochenderfer, M. W. M. Edwards, L. P. Espindle, J. K. Kuchar, and J. D. Griffith, "Airspace encounter models for estimating collision risk," *AIAA Journal of Guidance, Control, and Dynamics*, vol. 33, no. 2, pp. 487–499, 2010.
[13] M. J. Kochenderfer, J. P. Chryssanthacopoulos, and P. P. Radecki, "Robustness of optimized collision avoidance logic to modeling errors," in *Digital Avionics Systems Conference*, 2010.
[14] C. Tomlin, G. J. Pappas, and S. S. Sastry, "Conflict resolution for air traffic management: A study in multiagent hybrid systems," *IEEE Transactions on Automatic Control*, vol. 43, no. 4, pp. 509–21, 1998.
[15] S. M. LaValle, *Planning Algorithms*. Cambridge University Press, 2006.
[16] D. Bertsekas, *Dynamic Programming and Optimal Control*. Athena Scientific, 2005.
[17] J. Foster and F. B. Richards, "The Gibbs phenomenon for piecewise-linear approximation," *The American Mathematical Monthly*, vol. 98, no. 1, pp. 47–49, 1991.
[18] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT Press, 1998.
[19] S. Davies, "Multidimensional triangulation and interpolation for reinforcement learning," in *Advances in Neural Information Processing Systems*, 1996, pp. 1005–1011.
[20] J. A. Boyan and A. W. Moore, "Generalization in reinforcement learning: Safely approximating the value function," in *Advances in Neural Information Processing Systems*, 1995, pp. 369–376.