

Safety Verification of Neural Network Controlled Systems

A. Claviere¹ E. Asselin¹ C. Garion² C. Pagetti³

¹Collins Aerospace, France

²ISAE-SUPAERO, France

³ONERA, France

Scientific meeting DGA AID, November 3, 2020

Table of Contents

- 1 Introduction
- 2 System model
- 3 Problem definition
- 4 Reachability-based approach
- 5 Experiments
- 6 Conclusion and future work

Table of Contents

- 1 Introduction
- 2 System model
- 3 Problem definition
- 4 Reachability-based approach
- 5 Experiments
- 6 Conclusion and future work

The safety perspective

Neural network controlled system

The combination of a continuous-time dynamical system with a discrete-time neural network based controller.

What if such a system is considered as *safety critical*?



One has to show evidence that the system fulfills a set of *safety requirements*.

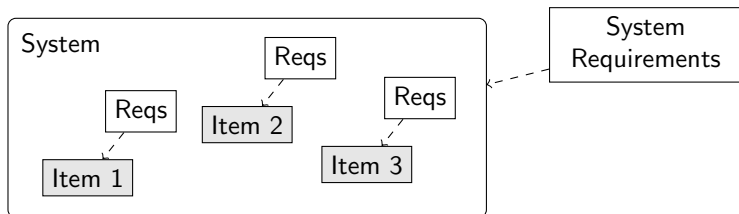
e.g., "A catastrophic failure shall occur with a probability less than 10^{-9} per hour of flight."

Safety: the classical approach

The system has to be developed in accordance with stringent *standards*.

They involve:

- Refinement of the system requirements at the *item* level
→ each item must be allocated a *correct, comprehensive* specification

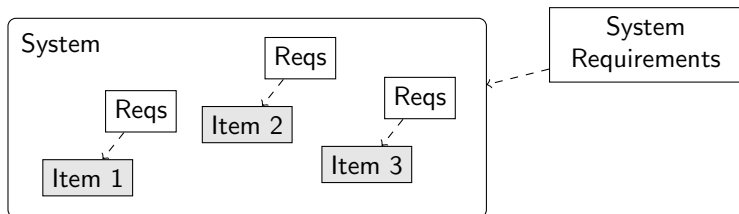


Safety: the classical approach

The system has to be developed in accordance with stringent *standards*.

They involve:

- Refinement of the system requirements at the *item* level
→ each item must be allocated a *correct, comprehensive* specification



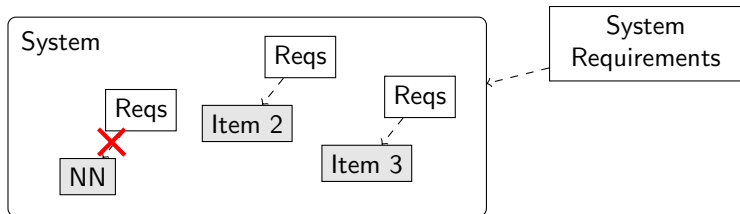
Each item must be developed in compliance with dedicated standards

→ each item must be shown to fulfill its specification

Safety: the case of NN controlled system

This approach is not applicable to neural network controlled systems:

- example data = *pointwise, non-comprehensive* specification
 → One cannot refine the system requirements at the network level



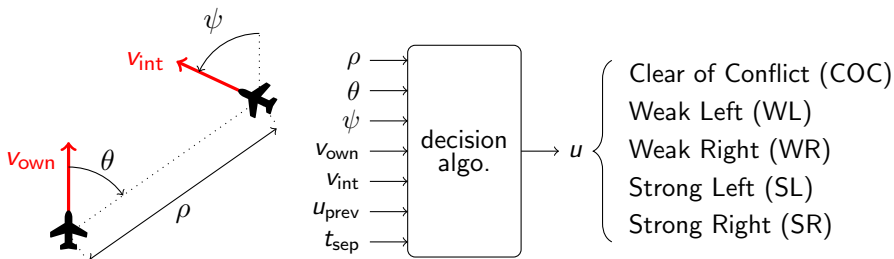
- The learning process does not guarantee the correctness of the network
 → It may be infeasible to show that a network fulfills its specification

Motivating example: the ACAS Xu controller

Two aircraft:

- the *ownship*, equipped with the ACAS Xu
- the *intruder*, equipped or not with the ACAS Xu

Objective: avoid a near mid-air collision between the two aircraft



Motivating example: the ACAS Xu controller

Two aircraft:

- the *ownership*, equipped with the ACAS Xu
- the *intruder*, equipped or not with the ACAS Xu

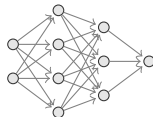
Objective: avoid a near mid-air collision between the two aircraft

Original design: lookup tables



2GB memory

Neural network approximation



2.4MB memory ($\times 0.001$)

relative runtime: $\times 0.97$

can be run on legacy avionics

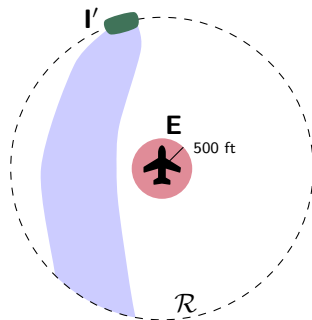
Motivating example: the ACAS Xu controller

Two aircraft:

- the *ownship*, equipped with the ACAS Xu
- the *intruder*, equipped or not with the ACAS Xu

Objective: avoid a near mid-air collision between the two aircraft

How to prove that the neural network based ACAS Xu is safe?



Safety: the case of NN controlled system

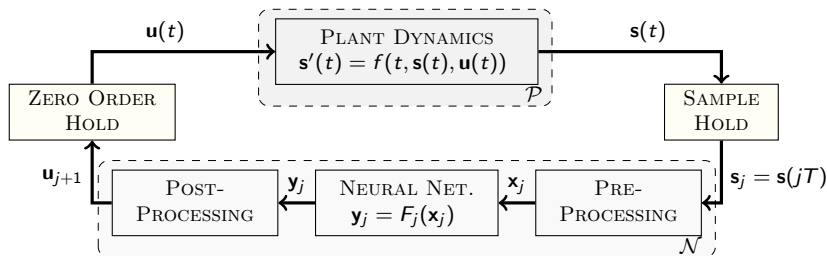
How to deal with these issues?

- Demonstrate safety without performing item-level refinement and analyses:
- construct a model of the overall system
 - perform a reachability analysis on this model

Table of Contents

- 1 Introduction
- 2 System model**
- 3 Problem definition
- 4 Reachability-based approach
- 5 Experiments
- 6 Conclusion and future work

Closed-loop system



Closed-loop system \mathcal{C} :

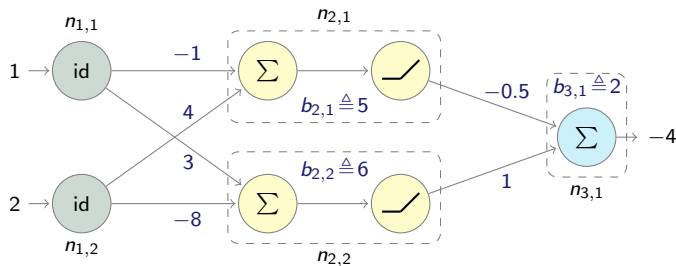
- a *continuous-time* plant \mathcal{P} , with state $\mathbf{s}(t) \in \mathbb{R}^l$
 - a *discrete-time* neural network based controller \mathcal{N} , with period T :
 - its j^{th} execution occurs in $[jT, (j+1)T[$
 - $\mathbf{u}(t) = \mathbf{u}_{j+1} \forall t \in [(j+1)T, (j+2)T[$
- \mathcal{N} is a *classifier* i.e., $\mathbf{u}_{j+1} \in \mathbf{U} = \{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(P)}\}$

Neural Network based controller

Uses a *collection* of ReLU networks $\mathbf{N} = \{N^{(1)}, \dots, N^{(D)}\}$
 \rightarrow only one network is executed at the j^{th} control step, depending on the previous command: $N_j = \lambda(\mathbf{u}_j)$

ReLU network

A (deterministic) function $F : \mathbb{R}^m \rightarrow \mathbb{R}^p$ that is the composition of affine transformations and non-linear ReLU units $\sigma : x \mapsto \max(0, x)$



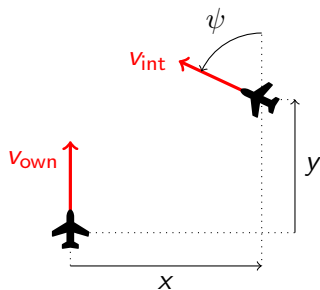
Plant - Example (ACAS Xu)

$$\mathcal{P} = \{\text{ownship, intruder}\}$$

$$\mathbf{s}(t) = (x(t) \ y(t) \ \psi(t) \ v_{\text{own}}(t) \ v_{\text{int}}(t))^T$$

simplified 2D kinematic model:

- the intruder has a uniform rectilinear displacement
- the ownship has a constant velocity



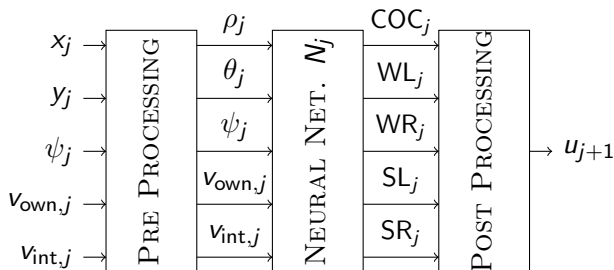
$$\left\{ \begin{array}{l} x'(t) = -v_{\text{int}}(t) \cdot \sin(\psi(t)) \\ y'(t) = v_{\text{int}}(t) \cdot \cos(\psi(t)) - v_{\text{own}}(t) \\ \psi'(t) = -u(t) \\ v'_{\text{own}}(t) = 0 \\ v'_{\text{int}}(t) = 0 \end{array} \right.$$

Neural Network based controller - Example (ACAS Xu)

Output = turn rate of ownship:

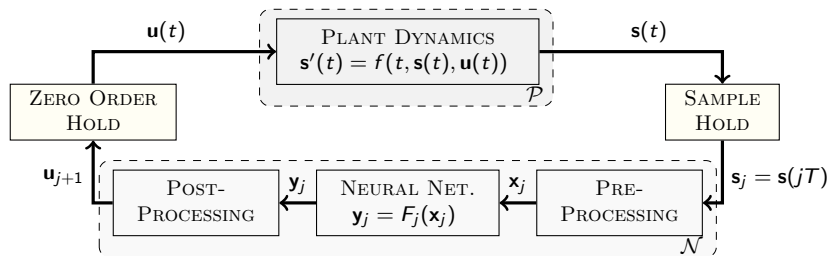
$$\mathbf{U} = \{0 \text{ deg/s}, 1.5 \text{ deg/s}, -1.5 \text{ deg/s}, 3 \text{ deg/s}, -3 \text{ deg/s}\}$$

Uses a collection of 5 ReLU networks $\mathbf{N} = \{N^{(1)}, \dots, N^{(5)}\}$



The post-processing is a argmin function

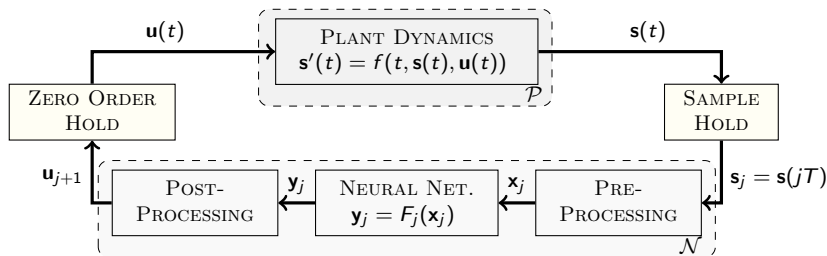
Closed-loop system



The state of the closed-loop \mathcal{C} is $\phi(t) = (\mathbf{s}(t), \mathbf{u}(t)) \in \mathbb{R}^l \times \mathbf{U}$

- $\mathbf{I} \subseteq \mathbb{R}^l \times \mathbf{U}$ is the set of the possible initial states
- $\mathbf{E} \subset \mathbb{R}^l \times \mathbf{U}$ is a set of *erroneous* states
- $\mathbf{T} \subset \mathbb{R}^l \times \mathbf{U}$ is a set of *target* states ($\mathbf{T} \cap \mathbf{E} = \emptyset$)

Closed-loop system



The state of the closed-loop \mathcal{C} is $\phi(t) = (\mathbf{s}(t), \mathbf{u}(t)) \in \mathbb{R}^l \times \mathbf{U}$

- $\mathbf{I} \subseteq \mathbb{R}^l \times \mathbf{U}$ is the set of the possible initial states
- $\mathbf{E} \subset \mathbb{R}^l \times \mathbf{U}$ is a set of *erroneous* states
- $\mathbf{T} \subset \mathbb{R}^l \times \mathbf{U}$ is a set of *target* states ($\mathbf{T} \cap \mathbf{E} = \emptyset$)

Deterministic behaviour: for a given initial state $\phi_0 \in \mathbf{I}$ and time horizon τ , there is a unique function ϕ_{ϕ_0} such that $\phi_{\phi_0}(t)$ is the state of \mathcal{C} at instant $t \leq \tau$

Closed-loop system - Example (ACAS Xu)

I: the ownship detects the intruder
($u = 0.0 \text{deg/s}$ *i.e.*, COC)

E: the intruder lies in the collision cylinder around ownship

T: the intruder is out of range

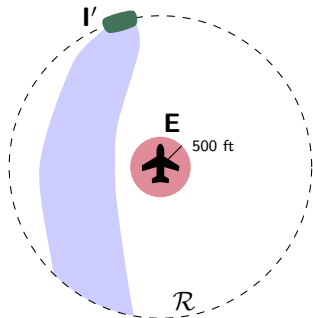


Table of Contents

- 1 Introduction
- 2 System model
- 3 Problem definition**
- 4 Reachability-based approach
- 5 Experiments
- 6 Conclusion and future work

Problem definition

Definition

The *safety verification problem* \mathcal{V} consists in deciding if:

$$\mathbf{R}_{[0,\tau]} \cap \mathbf{E} = \emptyset \quad (1)$$

where $\mathbf{R}_{[0,\tau]}$ are the reachable states over $[0, \tau]$

- \mathcal{V} is undecidable when the plant \mathcal{P} has a non-linear dynamics
- verifying pre/post-conditions on a network is a NP-hard problem
- the controller has a non-trivial logic

Problem definition

Definition

The *safety verification problem* \mathcal{V} consists in deciding if:

$$\mathbf{R}_{[0,\tau]} \cap \mathbf{E} = \emptyset \quad (1)$$

where $\mathbf{R}_{[0,\tau]}$ are the reachable states over $[0, \tau]$

Definition

The *safety verification problem* $\tilde{\mathcal{V}}$ consists in finding a set $\tilde{\mathbf{R}}_{[0,\tau]}$ satisfying $\tilde{\mathbf{R}}_{[0,\tau]} \supset \mathbf{R}_{[0,\tau]}$ and $\tilde{\mathbf{R}}_{[0,\tau]} \cap \mathbf{E} = \emptyset$

Table of Contents

- 1 Introduction
- 2 System model
- 3 Problem definition
- 4 Reachability-based approach**
- 5 Experiments
- 6 Conclusion and future work

Symbolic state and symbolic set

Definition

A *symbolic state* is a 2-tuple $([\mathbf{s}], \mathbf{u})$ wherein $[\mathbf{s}] \subset \mathbb{R}^l$ is a l -dimensional box *i.e.*, the cartesian product of l intervals, and $\mathbf{u} \in \mathbf{U}$. It symbolically represents the set $\{\phi(t) = (\mathbf{s}(t), \mathbf{u}(t)) \in \mathbb{R}^l \times \mathbf{U} \mid \mathbf{s}(t) \in [\mathbf{s}] \wedge \mathbf{u}(t) = \mathbf{u}\}$.

Definition

A *symbolic set* is a collection of symbolic states defined by $\tilde{\Phi} = \{([\mathbf{s}]_k, \mathbf{u}_k)\}_{1 \leq k \leq K}$ wherein $K \in \mathbb{N}$. It corresponds to the union of the sets represented by each $([\mathbf{s}]_k, \mathbf{u}_k)$.

Procedure

The procedure involves two types of sets:

- The symbolic set $\tilde{\mathbf{R}}_j$ approximates the reachable states at $t = jT$
- The symbolic set $\tilde{\mathbf{R}}_{[j]}$ approximates the reachable states for $t \in [jT, (j+1)T[$

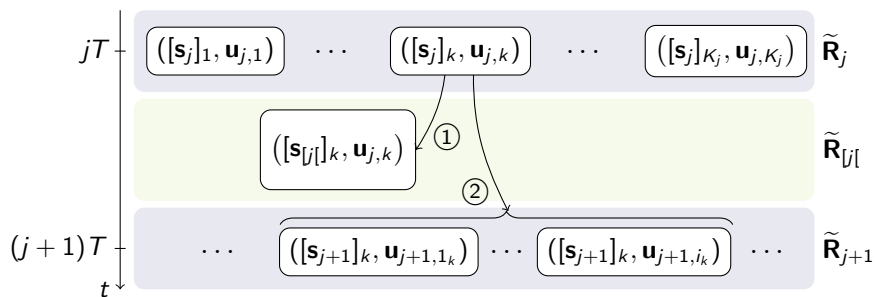
It works iteratively:

- starts with the symbolic set $\tilde{\mathbf{R}}_0 \supset \mathbf{I}$ enclosing the possible initial states.
- at control step j , builds $\tilde{\mathbf{R}}_{[j]}$ and $\tilde{\mathbf{R}}_{j+1}$ based on $\tilde{\mathbf{R}}_j$.

Finally, $\tilde{\mathbf{R}}_{[0,\tau]}$ is taken as the union of the $\tilde{\mathbf{R}}_{[j]}$.

Procedure

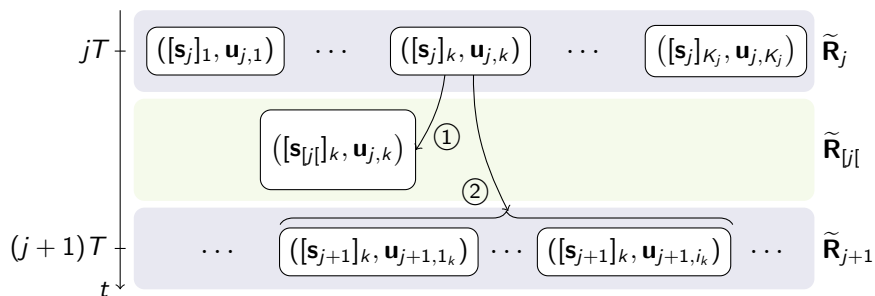
Approximation of the j^{th} control step:



- ① involves validated simulation (DynIBEX)
- ② involves both validated simulation (DynIBEX) and abstract interpretation (with a specialized solver ReluVal)

Procedure

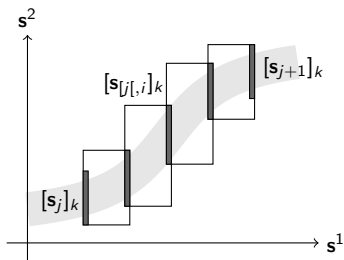
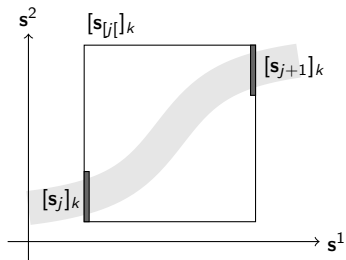
Approximation of the j^{th} control step:



To take account of a potential termination of \mathcal{C} ($\phi \in \mathbf{T}$), the symbolic states $([s_j]_k, \mathbf{u}_{j,k}) \subset \mathbf{T}$ are not further propagated

Optimizations

- Improving precision



- Improving time complexity

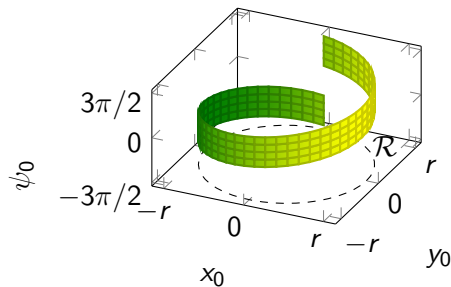
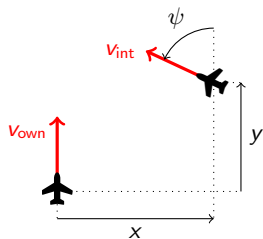
Avoid an exponential blow up of the number of symbolic states in $\tilde{\mathbf{R}}_j$
 → merge the “closest” symbolic states

Table of Contents

- 1 Introduction
- 2 System model
- 3 Problem definition
- 4 Reachability-based approach
- 5 Experiments**
- 6 Conclusion and future work

Experimental setup

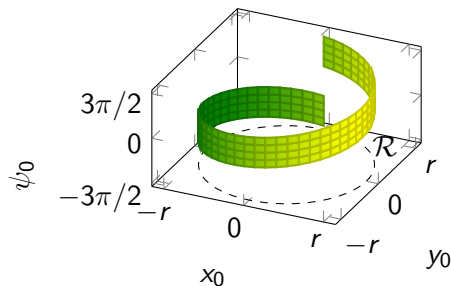
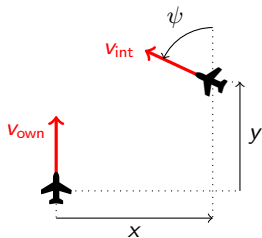
Partitioning: $\tilde{\mathbf{R}}_0 = \{([s_0]_k, 0.0 \text{ deg/s})\}_{1 \leq k \leq K_0}$ with $K_0 = 198,764$



- a single initial symbolic state $([s_0], 0.0 \text{ deg/s})$ approximating \mathbf{I} is unsafe
- the K_0 initial symbolic states composing $\tilde{\mathbf{R}}_0$ can be seen as K_0 *independent* verification problems (parallelization)
- the smaller the box $[s_0]_k$, the more precise the reachability analysis

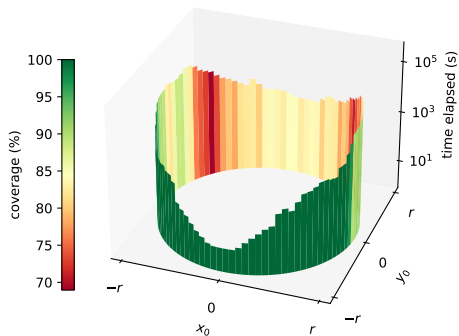
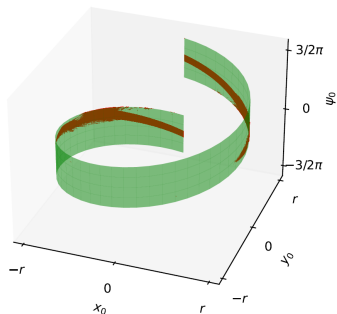
Experimental setup

Partitioning: $\tilde{\mathbf{R}}_0 = \{([s_0]_k, 0.0 \text{ deg/s})\}_{1 \leq k \leq K_0}$ with $K_0 = 198,764$



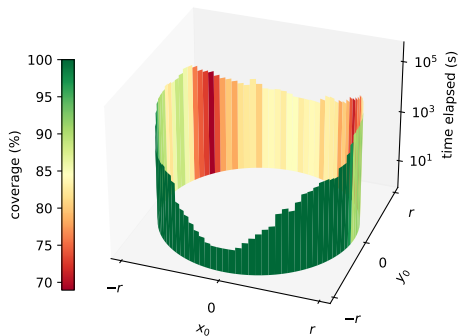
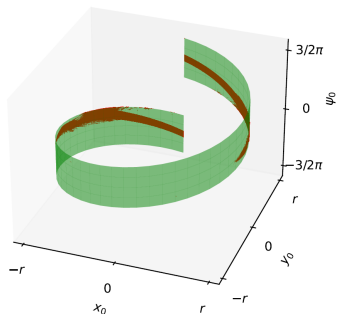
Split refinement: the initial symbolic states $([s_0]_k, 0.0 \text{ deg/s})$ for which the system cannot be proved safe are bisected \rightarrow new reachability analysis

Results



- coverage of 90.3%
- capability to identify the initial states for which the system could not be proved safe

Results



- most critical situations: the intruder is approaching from the left
- symmetry *w.r.t.* the $x_0 = 0$ axis: captures the symmetry of the collision avoidance problem

Table of Contents

- 1 Introduction
- 2 System model
- 3 Problem definition
- 4 Reachability-based approach
- 5 Experiments
- 6 Conclusion and future work**

Conclusion

- system-level approach for verifying the safety of neural network controlled systems
→ realistic model together with reachability analysis
- applicable to real-world systems
- provide valuable information from a practical point of view

Future work

- more efficient partitioning strategy (e.g., CFD)
- more efficient heuristics for splitting the initial symbolic states
- combine the approach with an efficient falsification strategy
- ACAS Xu: consider multiple UAVs, each one being equipped with a collision avoidance controller