

## Safe Exploration Algorithms for Reinforcement Learning Controllers

Mannucci, Tommaso; Van Kampen, Erik Jan; De Visser, Cornelis; Chu, Qiping

**DOI**

[10.1109/TNNLS.2017.2654539](https://doi.org/10.1109/TNNLS.2017.2654539)

**Publication date**

2018

**Document Version**

Accepted author manuscript

**Published in**

IEEE Transactions on Neural Networks and Learning Systems

**Citation (APA)**

Mannucci, T., Van Kampen, E. J., De Visser, C., & Chu, Q. (2018). Safe Exploration Algorithms for Reinforcement Learning Controllers. *IEEE Transactions on Neural Networks and Learning Systems*, 29(4), 1069-1081. <https://doi.org/10.1109/TNNLS.2017.2654539>

**Important note**

To cite this publication, please use the final published version (if applicable). Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

# Safe Exploration Algorithms for Reinforcement Learning Controllers

**Abstract**—Self-learning approaches, such as Reinforcement Learning, offer new possibilities for autonomous control of uncertain or time-varying systems. However, exploring an unknown environment under limited prediction capabilities is a challenge for a learning agent. If the environment is dangerous, free exploration can result in physical damage or in an otherwise unacceptable behavior. With respect to existing methods, the main contribution of this work is the definition of a new approach that does not require global safety functions, nor specific formulations of the dynamics or of the environment, but relies on interval estimation of the dynamics of the agent during the exploration phase, assuming a limited capability of the agent to perceive the presence of incoming fatal states. Two algorithms are presented with this approach. The first is the Safety Handling Exploration with Risk Perception Algorithm (SHERPA), which provides safety by individuating temporary safety functions, called backups. SHERPA is shown in a simulated, simplified quadrotor task, for which dangerous states are avoided. The second algorithm, denominated OptiSHERPA, can safely handle more dynamically complex systems for which SHERPA is not sufficient through the use of safety metrics. An application of OptiSHERPA is simulated on an aircraft altitude control task.

**Index Terms**—Reinforcement Learning, Safe Exploration, Adaptive Controllers, Model-free Control

## I. INTRODUCTION

IN engineering, classic control schemes such as PID still enjoy widespread use. This can be partially explained by the amount of effort needed to provide affordable yet efficient dynamic models of complex platforms. In the wake of this consideration, special attention in the control community has been dedicated to control schemes which require less precise knowledge of a model to achieve satisfactory performances. Robust control [1] constitutes an example of controller design developed to tolerate modeling error while guaranteeing a lower bound on performance. Adaptive control represents a promising field in developing new controllers with increased performance and reduced model dependency [2], [3]. A model-free option amidst adaptive control is Reinforcement Learning (RL).

Reinforcement learning is a knowledge based control scheme that mimics animal development [4]. At any given moment, an animal receives an array of internal and external stimuli that form its *situation*, with the *behavior* dictating the reaction to each of them. Correct reactions generate a positive chemical discharge that reinforces the behavior whereas unsuccessful ones lead to anguish that disproves it. This has an equivalent in Reinforcement Learning: stimuli constitute the *plant* or *system*; the animal is the *agent* following a temporary behavior, i.e. a *policy*. Lastly, the chemical reaction is represented by a numerical feedback called the *reward*.

From a theoretical point of view, RL has evolved in time to guarantee minimal level of performance. Selected algorithms

were proven to be Probably Approximately Correct (PAC) [5], and proofs of near-optimality and optimality for both discrete and continuous [6] applications were found. RL has also proved its worth in combination with Neural Networks (NN) in the development of Neuro-Dynamic Programming [7], [8], where the approximation power of NNs is utilized to efficiently represent both the value and the policy, thus reducing the “curse of dimensionality” [9]. In recent years, multiple successful applications [10]–[13], utilize RL as a framework for learning, in particular *model-free* RL [14]–[16].

In model-free RL there is no need of a model of the plant: learning starts with an *exploratory* phase, during which, rather than following the best policy based on current knowledge, new actions are tried in order to find the most rewarding actions and iteratively correcting its policy. The agent transitions gradually from exploring to *exploiting*, i.e. performing actions suggested by its best policy. The transition must be handled with care. Transitioning before the best policy is converged results in sub-optimal behavior due to lack of knowledge of the environment, but transitioning long after convergence unnecessarily delays the application of the best policy, thus affecting performance. This conflict between gathering and use of knowledge is called *exploration-exploitation dilemma*.

Whenever an agent attempts an inappropriate action, the consequent penalty acts as a negative reinforcement, and the wrong behavior is progressively discouraged, until it is no longer adopted. This permits to accommodate unwanted, unsafe actions within the canonical RL framework. However, consider an hostile environment where the consequences of wrong actions are not limited to “bad” performance, but include long term effects that can’t be compensated by more profitable exploitation later on. As one wrong action can result in unrecoverable effects, such an environment poses a *safety-exploration* dilemma, especially for a model-free approach. The goal of this work is to avoid such occurrences during learning, thus achieving *safe exploration* [17].

In the literature, a widely adopted approach to safety consists of assigning negative reward for undesired transitions, such that the most reliable policy maximises the minimal sum of reward in the presence of uncertainties and stochasticity. Safety is therefore embedded into policy performance. This *worst-case* or *minimax* [18] approach belongs to the *optimisation criterion* of safety [19]. Under this criterion, assuming a sufficiently large penalty for unsafe transitions, the optimal policy is also the safest. Methods for policy improvement with this approach have also been designed [20]. Alternatively, by assigning negative reward, the variance of the return can be taken into account by adopting *risk-sensitivity* approaches [21]–[24]. Several techniques exist to implement both minimax and risk-sensitive methods (see e.g. [25]); however, there

are limitations as far as exploration is considered. Including safety as part of the reward can generate a conflict between purely performance-based reward and safety-based reward if the penalty for unsafe transitions is not correctly assigned. Also, the optimisation criterion can be effective in preventing harmful situations, but requires previous knowledge of the probability of risk for the state-action space, which is in general the result of exploration itself. A different solution is to include safety in the *exploration process* itself. Garcia [19] differentiates three different approaches: “providing initial knowledge”, directing the learning in its initial stage towards more profitable and safer regions of the state space [26]; “deriving a policy from demonstrations” by means of *Learning from Demonstration* [27]; “providing teacher advice” by including an external teacher that can interrupt exploration and provide expert knowledge, or that the agent can consult when confronted with unexpected situations [28]. An alternative implementation of this solution is *risk-directed* exploration. With this approach, the agent’s choice of actions is aided by an appropriate *risk metric* [29] acting as an exploration bonus towards the safer regions of the environment.

Among algorithms that directly avoid unsafe transitions, [30] relies on an a-priori known safety function (acting as a go/no-go decision maker over feasible actions) and a fixed backup policy valid in all workspace. A similar approach is taken in [31], with the difference that the safety function is obtained through a “cautious simulator”. The simulator must correctly label unsafe states, but is allowed to mislabel safe states as unsafe: it is assumed that an experienced human operator can force the system into a mislabeled safe state. In [32] variable amount of perturbation is introduced in a given safe but inefficient baseline controller, so that discovery of new trajectories for task completion is possible, taking a certain amount of risk. These techniques share the need of a guaranteed safe controller, simulators or backup policy in order to prevent catastrophic exploration when facing critical decisions. Moldovan and Abbeel [33] define safety in terms of ergodicity of the exploration, and introduce an algorithm that relies on believes of the system, but not on a predefined baseline policy or safe controller.

This paper adopts a different strategy for safe exploration that does not rely on a-priori known safety functions, cautious simulators, or in an explicit ergodicity of the system. Instead, a temporary safety function is generated at each time in the form of *backups*. A backup is essentially an escape route: a control sequence that can bring the system in a close neighbourhood of a state that the agent already visited in the past. By resorting to backups, this strategy eliminates the need for the above prerequisites. This is particularly apt for RL tasks where the agent’s knowledge is very limited, which makes these prerequisites more demanding.

Two algorithms implementing the strategy are presented in this paper, and compared to similar preexisting ones. The first algorithm is the Safety Handling Exploration with Risk Perception Algorithm (SHERPA). Given a user-defined exploratory policy and an uncertain model of the system, SHERPA searches backups satisfying a *closeness condition*. If the search fails, SHERPA replaces the policy action with

a safer alternative, effectively acting as a “filter” with respect to the exploratory policy. SHERPA is effective in all those cases where backups can be found with reasonable ease. For more challenging tasks, a second version called OptiSHERPA is proposed which uses *metrics* to assess the safety of all available actions. Additionally, OptiSHERPA includes a dedicated evasion strategy, that relies on the current belief on the state space to minimize risks while at the same time avoiding impending dangers. While the “filtering” effect of both algorithms could be applied in diverse contexts, it is in RL that they prove more useful, due to the usual lack of information about the task being performed.

For the most general scenario, coverage of the state space, and convergence to the optimal policy cannot be guaranteed when using the proposed algorithms, which promote safety of exploration more than its efficiency (e.g. does not follow the “optimism in the face of uncertainty” [34] criteria). However, the formulation of the proposed algorithms is flexible enough to address these drawbacks without major modifications, e.g. by relaxing the constraints for the individuation of a backup.

Two applications of the algorithms are shown in the form of simulations. SHERPA is applied to a simulated quadrotor UAV exploring an indoor environment. OptiSHERPA is used to control the elevator deflection of a fighter aircraft exploring its flight envelope. Both agents have a stochastic exploratory policy: the goal of the algorithms is then to enforce safety during exploration.

## II. FUNDAMENTALS

This section will expand the problem of safe exploration. First, a brief overview of the fundamentals of RL will be provided. The motivation will be presented and the problem statement will be formalized. Finally, the assumptions of this method will be summarized.

### A. Problem statement

Classic RL is defined on the Markov Decision Process (MDP) framework. An MDP schematically represents a task for an agent in an environment, and consists of a tuple of five elements: state, action, transition, reward and discount. Set  $\mathcal{S}$  contains the states  $\mathbf{x}$  of the environment. Set  $\mathcal{A}$  is the set of *actions*  $\mathbf{u}$  that the agent can select.  $\mathcal{D}$  is a mapping  $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  that  $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{S}$  and  $\forall \mathbf{u} \in \mathcal{A}$  assigns to triplet  $(\mathbf{x}, \mathbf{u}, \mathbf{x}')$  the probability of the environment transitioning from state  $\mathbf{x}$  to state  $\mathbf{x}'$  given action  $\mathbf{u}$ . Function  $\mathcal{R} : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  maps all transitions to a scalar  $r$ , the *reward*, which indicates the immediate benefit of the transition with respect to accomplishing the task. A policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is a mapping that  $\forall \mathbf{x}$  assigns a probability of selecting action  $\mathbf{u}$ , describing the behavior of the agent. The goal of the agent is to learn an optimal policy  $\pi^*$  that maximizes the sum of expected reward  $J^\pi(\mathbf{x})$  subjected to a discount  $\gamma$  which privileges short term reward.  $J$  is the *value* of  $\mathbf{x}$  with  $\pi$ .

Various algorithms [35] exist to solve this problem, all of which rely on the concept of *exploration*. Exploration consists in performing different actions in different states, observing the consequent reward, and using this knowledge

to progressively define a policy  $\pi$  which maximizes the total expected reward. Under certain conditions, exploration can eventually yield policy  $\pi^*$  or its approximation [6].

A *fatal occurrence* is defined as an unacceptable condition for the agent; for example if the agent is harmed, e.g. a crash or a failure, or if it cannot proceed further in its task. Safe exploration consists in preventing fatal occurrences.

### B. Definitions and assumptions

In order to define the approach to safe exploration, this section will introduce multiple assumptions: a framework on how to represent fatal occurrences in physical agents; a limited prediction capability, denominated *risk perception*; and a model for the computation of uncertain dynamics. The dynamics represented by  $\mathcal{D}$  define the set  $\mathcal{T}$  of feasible transitions  $\tau = (\mathbf{x}, \mathbf{u}, \mathbf{x}'), \mathbf{x}, \mathbf{x}' \in \mathcal{S}, \mathbf{u} \in \mathcal{A}$ . Let  $\mathcal{T}_{\text{fat}}$  represent the subset of those transitions generating fatal occurrences. The following will be assumed:

*Assumption 1:* If  $\exists \bar{\tau} = (\bar{\mathbf{x}}, \bar{\mathbf{u}}, \bar{\mathbf{x}}') \in \mathcal{T}_{\text{fat}}$ , and  $\exists \mathbf{x} \in \mathcal{S}, \exists \mathbf{u} \in \mathcal{A}, \exists \tau' = (\mathbf{x}, \mathbf{u}, \mathbf{x}') \in \mathcal{T}$ , then  $\tau' \in \mathcal{T}_{\text{fat}}$

This assumption is not new to the literature (e.g. [36]) and can be considered an extension of the Markov property to fatal occurrences: these are inherently related to *fatal states* and not to actions. Any action leading to the same state or condition would result in the same fatal occurrence. This allows to define the *fatal state space* (FSS) as

*Definition 1 (Fatal state space):*

$$\text{FSS} = \{\mathbf{x}' | \forall \mathbf{x} \in \mathcal{S}, \forall \mathbf{u} \in \mathcal{A}, \tau = (\mathbf{x}, \mathbf{u}', \mathbf{x}') \in \mathcal{T}_{\text{fat}}\}$$

and analogously the *safe state space* (SSS) as

*Definition 2 (Safe state space):*

$$\text{SSS} = \mathcal{S} \setminus \text{FSS}.$$

Complex systems may present multiple modalities of fatal occurrences; these will be informally defined as *risks*. It is reasonable to assume that only a subset of state components will be involved for each risk. Given the state space  $\mathcal{S} \subseteq \mathbb{R}^n$ , define as the *restricted state space* (RSS) of a risk the space defined by those components involved in the risk. Consider for example an aircraft for which two risks are defined: hitting the ground and getting damaged by the aerodynamic forces. The ground risk can be related to the position of the agent with respect to the ground, whereas the aerodynamic risk can be related to wind speed, deflections, and attitudes. Different features of the state vector are involved for each risk. The first RSS is the space of all possible positions. The second RSS is the space of all possible combinations of wind speed, deflections, and attitudes. It will be assumed that the possible modes of fatal occurrence are known:

*Assumption 2:* For each risk, the agent has a-priori knowledge of the relative RSS.

Such a formulation arguably comes natural to a designer, since it expresses the risks involved in the exploration, and is also very high-level, so that the designer is not required to consider dangerous conditions but only those that are certainly fatal. Those elements of an RSS that are also fatal form the a-priori unknown *restricted fatal state space* (RFSS) (see Fig.1). Now that a structure for the fatal occurrences is defined, assumption 3 will introduce a mean of prevention of risks, providing closure to the problem definition.

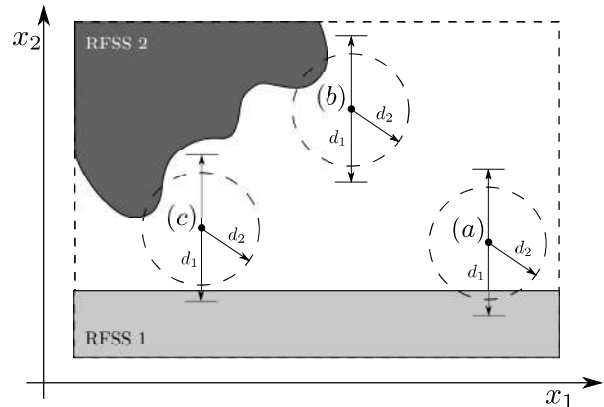


Fig. 1: An agent exploring a state space (bold dotted rectangle) with two different RFSS. The first cause of fatality is independent from the value of component  $x_1$ , and an Euclidean distance  $d_1$  is provided for the risk perception. A second cause of fatality is dependent on both components of the state so that an Euclidean distance  $d_2$  in the whole state space is given. In (a), RFSS1 is in reach; in (b), RFSS2 is in reach; in (c), both RFSSs are in reach. In all three cases the agent receives an identical warning.

*Assumption 3 (Risk perception):* Let  $\text{RSS}_k$  be the restricted state space of the  $k^{\text{th}}$  risk,  $k \in \{1, \dots, m\}$ . Let  $\mathbf{x}_k$  be the projection of state  $\mathbf{x} \in \mathcal{S}$  to  $\text{RSS}_k$ , and  $\mathbf{f}_k \in \text{RFSS}_k$  the nearest fatal state to  $\mathbf{x}_k$ . Then  $\forall k \in \{1, \dots, m\}, \exists \epsilon_k \in \mathbb{R}$  such that at each time-step, the value of the boolean function

$$\|\mathbf{x}_1 - \mathbf{f}_1\| \leq \epsilon_1 \vee \dots \vee \|\mathbf{x}_m - \mathbf{f}_m\| \leq \epsilon_m \quad (1)$$

is known, where  $\|\cdot\|$  is the canonical 2-norm.

The risk perception can be represented by a unknown mapping  $W : \mathcal{S} \rightarrow \{0, 1\}$ ,  $W(\mathbf{x})$  equal to 1 if (1) is true and 0 otherwise. Risk perception is a strong assumption; however, the following should be considered. When performing exploration in a dangerous environment, the agent cannot predict if the next state will be fatal or not without any form of knowledge of the FSS. Risk perception constitutes a plausible source of knowledge, arguably more valid than a-priori assumptions on fatal states distributions. Being a form of local, sensor based knowledge, risk perception allows to react to the insurgence of risks while retaining a sufficient level of abstraction and generality.

Consider the continuous-time version of an MDP. Indicating by  $\sigma(\mathbf{x}(t_0), \mathbf{u}(t), t)$  the state of the system at time  $t > t_0$  when applying generic action history  $\mathbf{u}(\mathbf{x}(t_0), t)$ :

$$\sigma(\mathbf{x}(t_0), \mathbf{u}(t), t) = \mathbf{x}(t_0) + \int_{t_0}^t \dot{\mathbf{x}}(\mathbf{x}, \mathbf{u}(\tau)) d\tau, \quad (2)$$

the set  $L$  of the lead-to-fatal (LTF) states is

*Definition 3 (Lead-to-fatal states):*

$$L = \{\mathbf{x} | \forall \mathbf{u}(\mathbf{x}(t_0), t), \exists t : \sigma(\mathbf{x}(t_0), \mathbf{u}(t), t) \in \text{FSS}\}. \quad (3)$$

LTF states are those that do not belong to the FSS, but evolve in the FSS with probability one. They can be seen as an equivalent in the generalized states space of the *inevitable collision states* introduced in [37], and as the *supercritical states* defined in [30] for a deterministic MDP.

In model-based RL, the controller has an adaptive model of the environment that allows off-line training and on-line system identification. Instead, in model-free RL (e.g. Q-learning [38]) there is no off-line training - the policy refinement is based on on-line training and no model is generated. An intermediate position between these two instances will be adopted in this work by introducing the following definition:

*Definition 4 (Bounding model):* Given dynamics  $\mathcal{D}$  yielding transition set  $\mathcal{T}$ , model  $\Delta(\mathbf{x}, \mathbf{u})$  is bounding for  $\mathcal{D}$  if and only if:

$$\forall \tau = \{\mathbf{x}, \mathbf{u}, \mathbf{x}'\} \in \mathcal{T}; \mathbf{x}' \in \Delta(\mathbf{x}, \mathbf{u}) \quad (4)$$

i.e. a model is bounding if it predicts *at least* all feasible transitions. Such a model is not used to perform off-line policy improvements but only to predict the boundaries of the dynamical evolution of the agent. It will be assumed that both real and bounding models are continuous, deterministic, and time-invariant.

### III. SHERPA

This section presents SHERPA. In order to properly explain the algorithm, *Interval Analysis* (IA) will be briefly introduced as a mean of obtaining bounding models [39]. The mathematical framework will then be expanded and refined as well. The algorithm will then be discussed in detail. Finally, the section will present an application.

#### A. Interval analysis

An interval  $[\mathbf{x}] = [\underline{\mathbf{x}}, \bar{\mathbf{x}}]$  is the set  $\{\mathbf{x} | x_i \leq \bar{x}_i, i \leq n\}$ ,  $\underline{\mathbf{x}} = (\underline{x}_1, \dots, \underline{x}_n)^T$ ,  $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_n)^T$ ,  $\underline{x}_i, \bar{x}_i \in \mathbb{R}$ . If  $\underline{\mathbf{x}} = -\bar{\mathbf{x}}$ ,  $[\mathbf{x}]$  is *symmetric*. In IA, the generic operator  $*$  applied to intervals yields all feasible solutions of the same operation applied to elements of the intervals:  $\forall \mathbf{x} \in [\mathbf{x}], \mathbf{y} \in [\mathbf{y}], \mathbf{x} * \mathbf{y} \in [\mathbf{x}] * [\mathbf{y}]$ . Consider a parametric formulation of  $\mathcal{D}$  as  $\mathcal{D}(\mathbf{x}, \mathbf{u}, \rho_1, \dots, \rho_k)$ , with parameters  $\rho_i \in [\rho_i]$  and  $[\rho_i]$  known confidence intervals. Following the above definition,  $[\mathcal{D}(\mathbf{x}, \mathbf{u})] = \{\mathcal{D}(\mathbf{x}, \mathbf{u}, \rho_1, \dots, \rho_k) | \rho_i \in [\rho_i]\}$  is the bounding model of  $\mathcal{D}$ . For each  $\mathbf{x}$  and  $\mathbf{u}$ ,  $[\mathcal{D}(\mathbf{x}, \mathbf{u})]$  is also an interval. This makes bounding models relatively simple and very economic in terms of data representation. However, it has the disadvantage of overestimating the *reachable set* of the states, due to an effect known as *dependency problem*, which propagates uncertainty in states and parameters.

#### B. Background and algorithm description

The final goal of the agent is to select at each time  $t$  a control  $\mathbf{u}$  that will keep the system safe, i.e. will avoid the FSS. Safety of a control can be defined as follows:

*Definition 5:* Control  $\mathbf{u}(t)$ ,  $t \in [t_1 t_2]$  is **safe** in state  $\mathbf{x}(t_1)$  if

$$\sigma(\mathbf{x}(t_1), \mathbf{u}, t) \in \text{SSS}, \forall t \in [t_1 t_2].$$

Safe control  $\mathbf{u}(t)$  is *feasible* if  $\forall t, \underline{\mathbf{u}} \leq \mathbf{u}(t) \leq \bar{\mathbf{u}}$ , where  $\underline{\mathbf{u}}$  and  $\bar{\mathbf{u}}$  represent boundaries on control, e.g. saturation. If this condition is strictly satisfied,  $\mathbf{u}(t)$  is said to be *strictly feasible*. However, multiple limitations such as inertia, control saturation and holonomic constraints give rise to LTF states, which must also be avoided during exploration, but are not

perceivable. For a generic trajectory  $\sigma(\mathbf{x}(t_0), \mathbf{u}, t)$ , the condition of (3) cannot be verified in practice since it would require an infinite time horizon. However, an exception occurs when a trajectory  $\sigma$  visits a state  $\mathbf{x}(\bar{t})$  multiple times. Formally, given a safe and feasible control  $\mathbf{u}(t)$  with associated trajectory  $\sigma(\mathbf{x}(t_0), \mathbf{u}, t)$ ,  $t \in [t_0 t_1]$  if

$$\exists t', t'' \in [t_0 t_1], t' \neq t'': \sigma(\mathbf{x}(t_0), \mathbf{u}, t') = \sigma(\mathbf{x}(t_0), \mathbf{u}, t'') \quad (5)$$

then there exists a safe and feasible control  $\mathbf{u}^*(t)$  for  $t \in [t_0 \infty)$ . It can be seen from (2) that (5) can be rewritten as  $\sigma(\mathbf{x}(t'), \mathbf{u}(t - (t' - t_0)), t') = \sigma(\mathbf{x}(t'), \mathbf{u}(t - (t' - t_0)), t'')$ . Then,  $\forall \Delta t \leq (t'' - t')$ ,  $\sigma(\mathbf{x}(t''), \mathbf{u}(t - (t'' - t') - (t' - t_0)), t'' + \Delta t) = \sigma(\mathbf{x}(t'), \mathbf{u}(t - (t' - t_0)), t' + \Delta t)$ . Since trajectory  $\sigma(\mathbf{x}(t'), \mathbf{u}(t), t' + \Delta t) \in \text{SSS}$ , a safe and feasible control exists at least for  $t \in [t_0 t'' + (t'' - t')]$ . The above can be repeated by replacing  $t'$  with  $t''$  and  $t''$  with  $t'' + (t'' - t')$ , guaranteeing a safe control for  $t \in [t_0 t'' + 2(t'' - t')]$ . The procedure can be repeated indefinitely, thus demonstrating the above. A particular occurrence of (5) consists in those cases where the system is in *equilibrium*, i.e.  $\dot{\mathbf{x}} = 0$ . Then, define an *ideal backup* from  $\mathbf{x}_0 = \mathbf{x}(t_0)$  to  $\mathbf{x}(\bar{t})$  as follows:

*Definition 6:* A feasible, safe control action  $\mathbf{u}_b(t)$ ,  $t \in [t_0 \bar{t}]$  is an **ideal backup** from  $\mathbf{x}_0$  to  $\mathbf{x}(\bar{t})$  if

$$\exists t < \bar{t}: \mathbf{x}(\bar{t}) = \mathbf{x}(t). \quad (6)$$

An ideal backup guarantees safety; however, (6) cannot be verified when taking into account the uncertainties of the bounding model, since it is not possible to exactly predict the next state. The following assumption will be made:

*Assumption 4:* Let  $\mathbf{u}(t)$  be a strictly feasible safe control for  $\mathbf{x}_0 = \mathbf{x}(t_0)$  and  $t \in [t_0 \bar{t}]$ . Then  $\forall \epsilon > 0$ ,  $\exists M \in \mathbb{R}$  and  $\mathbf{u}^*(t)$  a safe control for  $\mathbf{x}: \|\mathbf{x} - \mathbf{x}_0\| < \epsilon$  and  $t \in [t_0 \bar{t}]$  such that

$$\max_t \|\mathbf{u}(t) - \mathbf{u}^*(t)\| \leq M, \lim_{\epsilon \rightarrow 0} M = 0$$

i.e. in a neighborhood of a state with a strictly feasible safe control, a safe control for any state of the neighborhood can be found by altering the safe control by a finite amount that tends to zero as the neighborhood reduces in size<sup>1</sup>. A weaker definition of a *backup* can then be introduced:

*Definition 7:* A feasible, safe control action  $\mathbf{u}_b(t)$ ,  $t \in [t_0 \bar{t}]$  is a **backup** from  $\mathbf{x}_0$  to  $\mathbf{x}(\bar{t})$  with reach  $\epsilon$  if

$$\exists t < \bar{t}: \|\mathbf{x}(\bar{t}) - \mathbf{x}(t)\| \leq \epsilon \quad (7)$$

which in turn allows for the following theorem:

*Theorem 1:* Let  $\mathbf{x}_1 = \mathbf{x}(t_1) \in \text{SSS}$ , let  $\mathbf{u}_S(t)$  be a strictly feasible safe control for  $t \in [t_0 t_1]$ , and let  $\mathbf{u}_I(t)$  be a strictly feasible ideal backup from  $\mathbf{x}_1$  to  $\mathbf{x}_2 = \sigma(\mathbf{x}_1, \mathbf{u}_I, t_2)$ . Then  $\exists \mathbf{u}_b$  a backup from  $\mathbf{x}_1$  to  $\mathbf{x}_2$  with reach  $\epsilon > 0$  and a control  $\mathbf{u}^*(t)$  for  $t \geq t_2$  such that

$$\mathbf{u}(t) = \begin{cases} \mathbf{u}_b(t) & \text{if } t \leq t_2 \\ \mathbf{u}^*(t) & \text{otherwise} \end{cases}$$

is a feasible safe control for  $t \in [t_1 \infty)$ .

<sup>1</sup>Intuitively, local Lipschitz continuity of  $\frac{d\mathbf{x}}{dt}(\mathbf{x}, \mathbf{u})$  is required for Assumption 4 to hold. However, the remainder of the paper will consider Assumption 4 as given.

*Proof:* the existence of the ideal backup  $\mathbf{u}_I(t)$  means that the iterative control action  $\mathbf{u}_{iter}(t) = (\mathbf{u}_I(t), \mathbf{u}_S(t), \mathbf{u}_I(t), \dots)$  is a strictly feasible safe control. The existence of  $\mathbf{u}_I(t)$  means that  $\forall \epsilon$  there exists a backup. Assumption 4 guarantees that there exists a safe control  $\mathbf{u}(t)$  such that  $\forall t \|\mathbf{u}(t) - \mathbf{u}_{iter}(t)\| \leq M$  after the application of the backup. Since  $\mathbf{u}_{iter}(t)$  is strictly feasible, and since, for  $\epsilon \rightarrow 0$ ,  $M \rightarrow 0$  then  $\exists \epsilon$  s.t.  $\mathbf{u}(t)$  is also feasible.  $\square$

Theorem 1 is of limited practical use - since it relies on the existence of  $\mathbf{u}_I(t)$  and does not specify the actual value of  $\epsilon$  - but formalizes the observation contained in (5) in the presence of uncertainty. The equilibrium condition can be formalized as follows:

*Definition 8:* State  $\mathbf{x}_E \in \text{SSS}$  is an *equilibrium point* for the system if  $\exists \mathbf{u}_E = \mathbf{u}(t_E)$  s.t.  $d\mathbf{x}/dt(\mathbf{x}_E, \mathbf{u}_E) = \mathbf{0}$ .

The definition of backup can be revised as

*Definition 9:* A feasible, safe control action  $\mathbf{u}_b(t)$ ,  $t \in [t_0, \bar{t}]$  is a **backup** from  $\mathbf{x}_0$  to  $\mathbf{x}(\bar{t})$  with reach  $\epsilon$  if  $\exists t < \bar{t} : \|\mathbf{x}(\bar{t}) - \mathbf{x}(t)\| \leq \epsilon$ , in which case is said to have *reach*  $\epsilon$ , or if there exists an equilibrium point  $\mathbf{x}_E$  s.t.  $\|\mathbf{x}(t_2) - \mathbf{x}_E\| \leq \delta$ .

### C. Closeness condition

The goal of SHERPA is to find a feasible backup with a sufficiently small reach, i.e. for which the system reaches a “close” neighborhood of a previously visited state; therefore, a *closeness condition* must heuristically be introduced to implement the method. Once again, IA offers a simple way of interpreting closeness between two points. Before proceeding further, it is convenient to slightly alter the definition of backup to accommodate the use of intervals. Let the bounding model  $\Delta$  of  $\mathcal{D}$  be its IA extension  $[\mathcal{D}]$  with time-discrete formulation, so that  $[\mathbf{x}_{k+1}] = \Delta([\mathbf{x}_k], \mathbf{u}_k)$ , and let  $[\epsilon]$  and  $[\delta]$  be two  $n$ -dimensional symmetric intervals. Then, a backup can be reformulated as follows:

*Definition 10:* Let then  $\{\mathbf{x}_k, \dots, \mathbf{x}_{k+m}\}$  be the trajectory generated by control  $\mathbf{u}_b = \{\mathbf{u}_k, \dots, \mathbf{u}_{k+m}\}$ . Then  $\mathbf{u}_b$  is a **backup** in interval form for  $[\mathbf{x}_k]$  if and only if:

$$\forall i \in \{1, 2, \dots, m\}, [\mathbf{x}_{k+i}] \subset \text{SSS} \quad (8)$$

and :

$$\forall \mathbf{x}_{k+m} \in [\mathbf{x}_{k+m}], \exists p \leq k : (\mathbf{x}_{k+m} - \mathbf{x}_p) \in [\epsilon] \quad (9)$$

in which case the backup has reaching interval  $[\epsilon]$ , or if it exists an equilibrium point  $\mathbf{x}_E$  s.t. :

$$\forall \mathbf{x}_{k+m} \in [\mathbf{x}_{k+m}], (\mathbf{x}_{k+m} - \mathbf{x}_E) \in [\delta] \quad (10)$$

in which case the backup has reaching interval  $[\delta]$ .

This definition replaces the reach with intervals  $[\epsilon]$  and  $[\delta]$ . Given two symmetric threshold intervals  $[\epsilon]_{thr}$  and  $[\delta]_{thr}$ , a backup is effective if either  $[\epsilon] \subset [\epsilon]_{thr}$  or  $[\delta] \subset [\delta]_{thr}$ . Fig. 2 shows an example of backup evaluation. Length  $m$  should be chosen so that the size of intervals  $[\mathbf{x}_{k+m}]$  is comparable to risk perception ranges  $\epsilon$ .

The choice of the threshold intervals is empirical in nature; however  $[\epsilon]_{thr}$  should depend on the magnitude of the control.

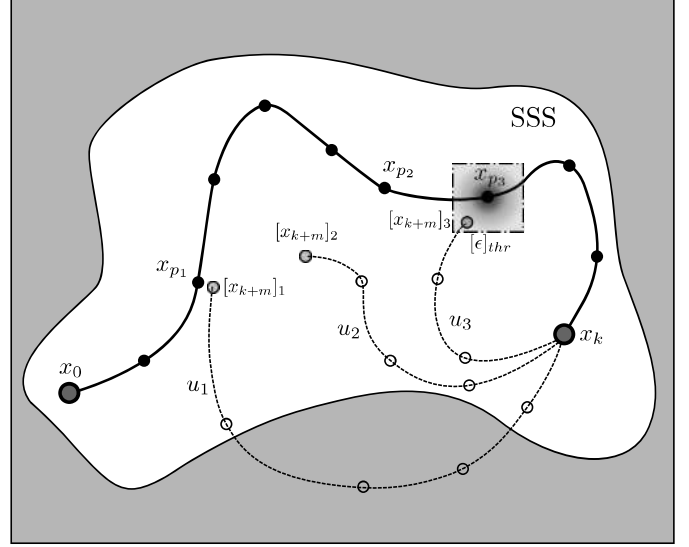


Fig. 2: An example of backup evaluation. Starting from position  $\mathbf{x}(0) = \mathbf{x}(t_0)$ , the agent has followed its trajectory to current state  $\mathbf{x}_k$ . The three dashed lines represent the means of three proposed interval backups from  $\mathbf{x}_k$ . The proposed backup generated by control  $\mathbf{u}_1$  is discarded since it violates the known safe state space SSS. The proposed backup generated by  $\mathbf{u}_2$  is also discarded since  $[\mathbf{x}_{k+m}]_2$  does not satisfy the closeness condition with nearest explored state  $\mathbf{x}_{p2}$ . Control  $\mathbf{u}_3$  generates a feasible backup, since the associated trajectory lies entirely in the SSS, and since  $[\mathbf{x}_{k+m}]_3$  respects the closeness condition.

Given assumption 4,  $\lim_{\epsilon \rightarrow 0} \max_k \|\mathbf{u}(t_k) - \mathbf{u}^*(t_k)\| = 0$ , so that the smaller the threshold  $[\epsilon]_{thr}$ , the more likely the feasibility of control  $\mathbf{u}(t_k)$ . To account for this, an initial threshold interval  $[\epsilon]_{max}$  is assigned to each new state  $\mathbf{x}_k$ . At each successive time step, all the threshold intervals are shrunk by a factor  $\eta$ :

$$\eta = \left( \prod_{i=1}^m \eta_i \right)^{\frac{1}{m}}; \quad \eta_i = \left( 1 - \frac{|u_i(t_k) - \frac{1}{2} \cdot (\underline{u}_i + \bar{u}_i)|}{\frac{1}{2} \cdot (\bar{u}_i - \underline{u}_i)} \right)^{\lambda} \quad (11)$$

with positive exponent  $\lambda < 1$ . The more the control action nears the boundaries  $\underline{u}$  and  $\bar{u}$ , the smaller the  $\eta$ , and the more the shrinkage of  $[\epsilon]_{thr}$ . Exponent  $\lambda$  also regulates shrinkage: the smaller  $\lambda$ , the less previous explored states are penalized with respect to recent ones. Values of  $\lambda$  and of  $[\epsilon]_{max}$  are closely related; e.g. a small  $\lambda$  should be paired with a small  $[\epsilon]_{max}$  and a high  $\lambda$  with an increased  $[\epsilon]_{max}$ .

### D. SHERPA

In this section, the SHERPA algorithm (Fig. 3) will be discussed in detail. At the start of the exploration, the agent is in state  $\mathbf{x}_0$ . It is assumed that  $W(\mathbf{x}_0) = 0$ . For each proposed action  $\mathbf{u}_p$  generated by policy  $\pi$ , SHERPA computes the predicted interval  $[\mathbf{x}_p] = \Delta(\mathbf{x}_0, \mathbf{u}_p)$ . If  $[\mathbf{x}_p] \notin \text{SSS}(t_0)$ , a different action is proposed, or the agent might encounter a fatal state. Otherwise, SHERPA checks the existence of a backup for  $[\mathbf{x}_p]$ : a random finite control sequence  $\mathbf{u}_b$  is generated, and for each step, its reach is checked as per Definition 10. If no backup is found after a set number of

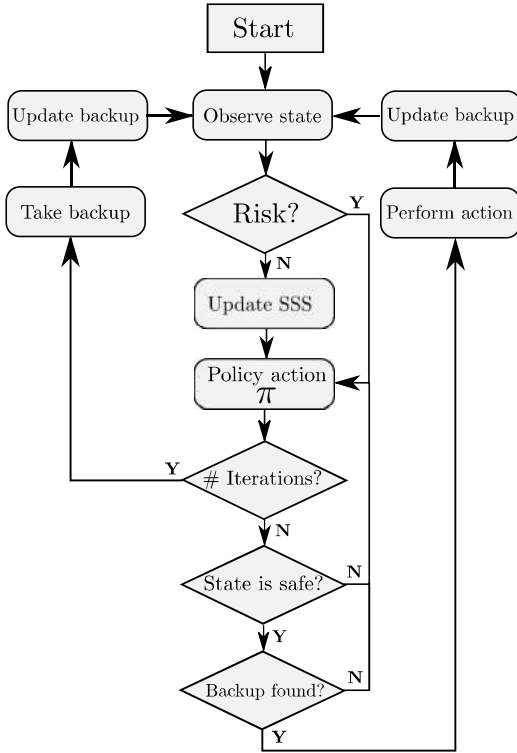


Fig. 3: Flowchart summarizing the procedure of SHERPA. The policy action is checked for safety and for backups; if the checks succeed, the action is taken and the new backup is stored in place of the previous. If the iteration limit is reached without success, SHERPA recollects the previously stored backup.

iterations, the current action is discarded, and a different action is proposed. In the new state  $x_1$ , if  $W(x_1) = 0$ , then the SSS is augmented with all the states currently in reach of the risk perception. The procedure is iterated during the whole exploration. In theory, the agent has always a feasible control at its disposal, eventually resorting to a backup, i.e. an iterative control, or a dynamical equilibrium. If all proposed actions fail the previous checks, SHERPA will force the agent to adopt the current backup. During backup execution, it will proceed to individuate a safe control and a backup for the expected arrival interval. The complexity of the algorithm depends on the number of visited states  $k$ , the number of available actions  $n_a$ , the maximum amount of iterations for action selections  $a_{max}$  and for backup checking  $b_{max}$ , and is equal to  $O(a_{max} \cdot \max(\min(b_{max}, n_a^m), kn^k))$ .

Other methods have proven their worth in preventing a controlled system from reaching a set of unsafe states or configurations. Potential fields [40] were successfully adapted for robotic navigation and path planning with lack of environmental knowledge, e.g. uncertainty in observed obstacles [41], in the number of obstacles [42] and even in the entire environment [43]. However, these methods are not equally adaptable when the uncertainty extends to the model dynamics, and when holonomic and control constraints are considered, as in this paper. An application of this is provided in the following as comparison with SHERPA. Lyapunov Barrier Functions [44], Control Barrier Functions [45] and Control

Lyapunov-Barrier Functions [46] can provide control that is guaranteed to be safe. However, the underlying assumptions on system structure, or on the barrier functions, can limit their applicability, especially when considering model uncertainty.

### E. Quadrotor task

In this section, a task is simulated for a two-dimensional quadrotor flying in a room. The underlying RL agent adopts a fully random policy to explore. The goal is for SHERPA to prevent collision and thus enable safe state exploration. The state vector is  $(x, \dot{x}, z, \dot{z}, \theta)^T$ , where  $x, z$  and  $\theta$  are the horizontal and vertical coordinate of the quadrotor, and the pitch angle. The bounding model<sup>2</sup> is

$$\begin{aligned} \ddot{x} &= -\frac{[c]}{[m]} |\sin \theta| |\dot{x}| \dot{x} + \frac{T}{[m]} \sin \theta; \quad \dot{\theta} = [k] I_\tau \\ \ddot{z} &= -\frac{[c]}{[m]} |\cos \theta| |\dot{z}| \dot{z} + \frac{T}{[m]} \cos \theta - g \end{aligned} \quad (12)$$

where  $c$  and  $[c]$  is a damping coefficient;  $m$  and  $[m]$  is the mass of the quadrotor;  $T$  is the total thrust generated by the rotors;  $k$  and  $[k]$  is an efficiency factor for torque impulse  $I_\tau$ ;  $g$  is the gravitational acceleration. Intervals  $[c]$ ,  $[m]$  and  $[k]$  are

$$\begin{aligned} [m] &= [0.344, 0.516] \text{ kg}; \quad [c] = [0.72, 1.08] \times 10^{-5} \frac{\text{kg}}{\text{m}}; \\ [k] &= [0.8, 1.2] \end{aligned}$$

with control action bounded as  $T \in [0, 12.7]$  N and  $I_\tau \in [-300^\circ, 300^\circ] \text{ s}^{-1}$ . The model of the quadrotor used for the computation of the real dynamics is obtained from the bounding model by randomly selecting  $c$ ,  $m$  and  $k$  from their intervals.

The quadrotor initially hovers at the middle of a room three meters wide and two meters high. Colliding with the walls or with the floor/ceiling is a fatal occurrence, and the quadrotor perceives risk half a meter from the ceiling or floor, or at one meter from a wall. The reach of the backup depends on  $[\epsilon^{max}] = ([-\frac{1}{2}, \frac{1}{2}]m; [-\frac{1}{2}, \frac{1}{2}] \frac{m}{s}; [-\frac{1}{2}, \frac{1}{2}]m; [-\frac{1}{2}, \frac{1}{2}] \frac{m}{s}; [-\frac{\pi}{6}, \frac{\pi}{6}])$ . During exploration, control is bounded between 10% and 90% of the original intervals:  $T \in [1.3, 10.3]$  N and  $I_\tau \in [-270^\circ, 270^\circ] \text{ s}^{-1}$ : this prevents factor  $\eta$  of (11) to vanish. Coefficient  $\lambda$  is chosen as 0.5. As for  $[\delta]$ , taking into consideration the dynamics in (12), it is evident that the position  $(x, y)$  of the quadrotor itself does not influence the value of the derivatives. Therefore the interval  $[\delta_{thr}] = ([-\infty, \infty]; [-1, 1] \frac{m}{s}; [-\infty, \infty]; [-1, 1] \frac{m}{s}; [-45, 45]^\circ)$  with equilibrium point  $\dot{x} = 0, \dot{z} = 0, \theta = 0$  is selected. SHERPA evaluates up to 10 actions per time step, proposed by the random policy. A maximum of 40 backup evaluations are performed, for a maximum of 400 iterations per time step. Initially, since the agent does not perceive any risk, the SSS( $t_0$ ) coincides with

$$([ -1, 1]m; [ -\infty, \infty]; [ -\frac{1}{2}, \frac{1}{2}]m; [ -\infty, \infty]; [ -\infty, \infty])$$

<sup>2</sup>It can be shown that, due to local Lipschitz continuity of (12), Assumption 4 is valid for any realization of the dynamics.

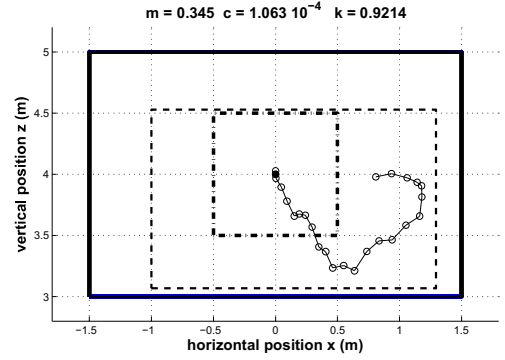
and it is updated whenever  $W = 0$  as

$$\text{SSS} \cup \{\mathcal{S}(x', z', \dot{x}, \dot{z}, \theta) \mid |x(t) - x'| < 1 \vee |z(t) - z'| < \frac{1}{2}\}.$$

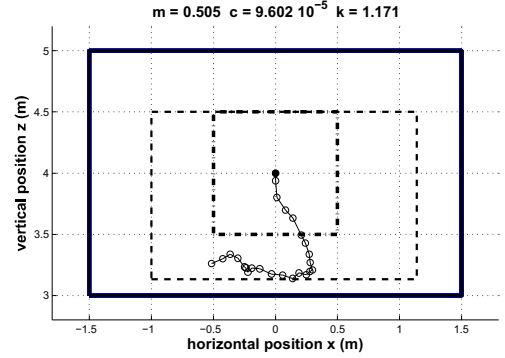
A series of 4 runs of SHERPA is shown in Fig. 4. The values of  $c$ ,  $m$  and  $k$  are indicated in the captions. The dots indicate the position of the quadrotor at different time steps, starting in position (0,4). The solid rectangle represents the fatal states. The dash-dotted line represents the contour of the region where no risk is individuated ( $W(x, y) = 0$ ) and the SSS is therefore updated. The dashed line represents the known SSS at the end of the run. In 3 out of 4 cases, the quadrotor reaches the boundaries of the SSS, but manages to safely divert its trajectory (Fig. 4.a and 4.b) or to wait in position until another path has been found (Fig. 4.c).

In general, the agent chooses actions that result in safe trajectories, even with a significant excursion in the value of parameters. Nonetheless, the random nature of the underlying policy is still noticeable. If a different task were to be proposed (e.g. waypoint navigation), SHERPA could still be utilized to evaluate actions suggested by a non exploratory, goal-oriented policy.

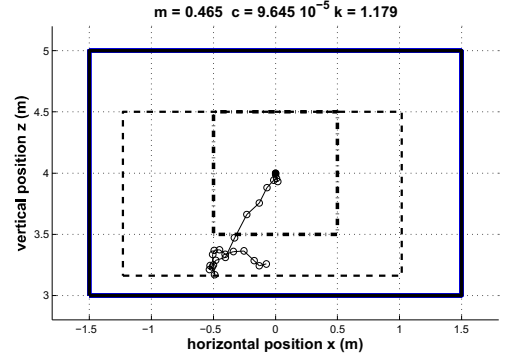
In order to highlight the novelty of the proposed approach, Fig. 5 and 6 show an application of a different method, specifically the Basic Potential Field (BPF) method [47] for collision avoidance. The BPF differs from other potential field methods such as [40] in that it generates a finite potential accounting for velocities and maximum decelerations. In this application, the BPF repels the quadrotor from exiting the known SSS. The rate at which the field reaches its maximum intensity is given by its gain  $G$  and by the worst-case maximum deceleration as given by the uncertain model. The quadrotor in Fig.5 is in an almost nominal condition of mass and torque effectiveness (as in Fig. 4d). It can be seen how a BPF with gain  $G = 0.5$  prevents collisions. The trajectory is centered around the middle of the known SSS, where the combined repulsive field, averaged between the different velocities, has a minimum. This differs from the proposed approach, in which the outer region of the SSS is not penalized with respect to its interior (see e.g. Fig.4a). Fig.6 shows the resulting trajectory when applying the previous BPF to a quadrotor with increased mass and torque effectiveness, as in Fig.4b. The figure shows how the quadrotor flies predominantly in the lower half of the room, as a result of the increased mass of the quadrotor. The trajectory violates the SSS, and even ends with a collision. This is due to the following. First, the quadrotor is forced to continuously change its attitude so that resultant of the forces complies with the direction indicated by the BPF. This essentially turns the BPF into a reactive collision avoidance method. As a result of this, the effectiveness of the method depends on the gain  $G$ , whose value determines the reactivity of the field, as shown by the two different conditions of Fig.5 and 6. If  $G$  is too low, the BPF is not sufficiently reactive; however, too high a gain can cause instability and unwanted oscillations. This selection of  $G$  is especially critical when considering systems with model uncertainties. This constitutes another difference with SHERPA, which takes into account all possible realizations of the model due to its interval formulation.



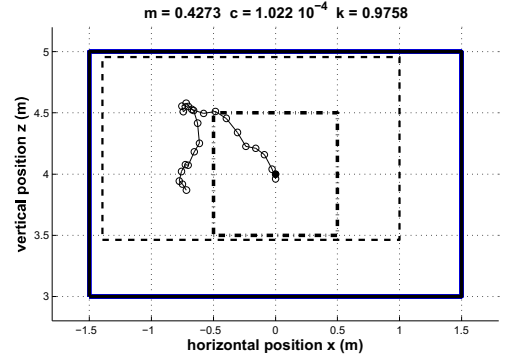
(a) Run with  $m = 345\text{g}$ ,  $c = 1.06 \times 10^{-4} \frac{\text{kg}}{\text{m}}$ ,  $k = 0.92$



(b) Run with  $m = 505\text{g}$ ,  $c = 9.60 \times 10^{-5} \frac{\text{kg}}{\text{m}}$ ,  $k = 1.17$



(c) Run with  $m = 465\text{g}$ ,  $c = 9.64 \times 10^{-5} \frac{\text{kg}}{\text{m}}$ ,  $k = 1.18$



(d) Run with  $m = 427\text{g}$ ,  $c = 1.02 \times 10^{-4} \frac{\text{kg}}{\text{m}}$ ,  $k = 0.98$

Fig. 4: The dots represent the quadrotor position in time. The solid rectangle represents the fatal states that the agent must avoid. The dot-and-dashed line delimits the states where the agent does not perceive risk. The dashed line represents the known SSS at the end of each run.



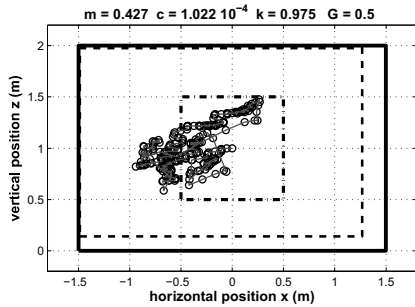


Fig. 5: The trajectory obtained when using the BPF to control a quadrotor in an almost nominal condition of mass and torque efficiency.

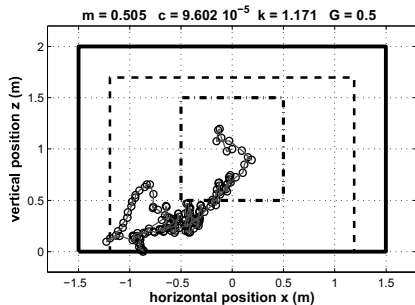


Fig. 6: The trajectory obtained when using the BPF to control a quadrotor with increased mass and torque efficiency. The quadrotor flies predominantly in the lower half of the room before colliding with the floor.

## IV. OPTISHERPA

### A. Motivation and algorithm description

The previous section introduced SHERPA, which relies on an a-priori defined interval to define closeness. This approach was shown to succeed in the quadrotor application. However, this approach has two drawbacks. First, if the system is not easily controllable, predicted arrival intervals might fail in satisfying the closeness condition. For these systems a safe backup could involve a complex trajectory which would require a considerable amount of time to be followed. The trajectory might also reach a portion of the state space outside of the risk-perception range of SHERPA, thus requiring a certain degree of “global” knowledge to be ensured as safe. The confidence in state prediction and the current knowledge of the SSS (provided by the risk-perception) might not be sufficient to identify such a trajectory as a backup. This would lead to situations where SHERPA is either unable to find a suitable control, thus exhausting its backup and then rejecting all further actions, or is forced to keep the system on hold in an equilibrium point.

A second drawback of the closeness condition is that, acting as a yes-or-no filter, it does not distinguish between control actions that are almost satisfactory and those that are completely unacceptable. After resorting to a backup, and while trying to find a new one from the arrival condition, SHERPA has only a limited amount of iterations to find a new backup. In the event that this search is unsuccessful, SHERPA is forced to take a possibly unsafe action.

For the above reasons, a second version of SHERPA named OptiSHERPA is presented. It differs from SHERPA in two ways. First, OptiSHERPA introduces metrics for the selection of actions. A finite set of actions is evaluated at each time-step together with a feasible backup, and the best action is performed. The introduction of metrics reduces the burden of on-line application by allowing OptiSHERPA to rank its options and to take an informed decision if the available amount of on-line iterations is depleted. A second difference lies in the strategy itself. With SHERPA, the generation of backups would keep the system safe, providing a possible escape route at every time-step. With OptiSHERPA, when the agent does not perceive danger, a *distance metric* is implemented with the goal of preventing potentially unsafe behavior of the system, effectively constraining the dynamics similar to SHERPA. When danger is detected, the agent examines its current belief in the state space and actively avoids the regions of the state space that are least safe. This is done via an *evasion metric*.

The remainder of this section is as follows. First, the metrics will be discussed in detail. Second, the algorithm will be illustrated. Third, a simulated application to an elevator control task will be shown.

### B. Metrics

1) *Distance metric*: This metric allows to classify intervals based on distance so that, during the evaluation of backups, the one with “closest” reach can be selected. Indicating by  $\odot$  the Schur product,

$$d(\mathbf{x}, \mathbf{I}) = \left\| \left( \mathbf{x} - \frac{\bar{\mathbf{I}} + \underline{\mathbf{I}}}{2} \right) \odot \mathbf{v}_r \right\| + \rho \cdot \left\| \left( \frac{\bar{\mathbf{I}} - \underline{\mathbf{I}}}{2} \right) \odot \mathbf{v}_r \right\| \quad (13)$$

is a “distance”<sup>3</sup> between  $\mathbf{x}$  and the center of the interval, rescaled by  $\mathbf{v}_r \in \mathbb{R}^{n+}$ , plus a term proportional to the interval width weighted by a positive parameter  $\rho < 1$ . This term allows to include the uncertainty in the interval as a penalizing factor. A lower value of  $\rho$  privileges intervals whose center are nearer regardless of their width; conversely a higher value penalizes intervals whose center is nearer but whose elements are more dispersed. Fig. 7 shows this transition. Vector  $\mathbf{v}_r$  should be chosen so that, for any two excursions in state  $\Delta \mathbf{x}^I$  and  $\Delta \mathbf{x}^{II}$ , it is  $\|\Delta \mathbf{x}^I \odot \mathbf{v}_r\| > \|\Delta \mathbf{x}^{II} \odot \mathbf{v}_r\|$  iff excursion  $\Delta \mathbf{x}^I$  affects safety and controllability of the system more than excursion  $\Delta \mathbf{x}^{II}$ . Thus, to components  $x_j$  related to risks should correspond an adequately big vector component  $v_j$ . A second function of  $\mathbf{v}_r$  is to normalize distances in  $\mathcal{S}$ , since components  $x_i$  might have different units of measure. The magnitude of the control can also be accounted for by a similar method as the one illustrated in (11). At each time step compute  $\eta$  so that the metric distance between  $\mathbf{x}$  and interval  $\mathbf{I}(t_i)$  at time  $t_i$  is:

$$d_\eta(\mathbf{x}, \mathbf{I}(t_i), t_j) = d(\mathbf{x}, \mathbf{I}(t_i)) / \prod_{k=i}^j \eta(t_k) \quad (14)$$

<sup>3</sup>Eq. (13) is not rigorously a distance, as it is not defined on  $\mathbb{R}^n \times \mathbb{R}^n$  but on  $\mathbb{R}^n \times \mathbb{I}^n$ , where  $\mathbb{I}^n$  is the set of  $n$ -dimensional intervals. However, if (13) is restricted to the subset of “crisp” intervals, it is analogue to the Euclidean distance in  $\mathbb{R}^n$ , hence the denomination.

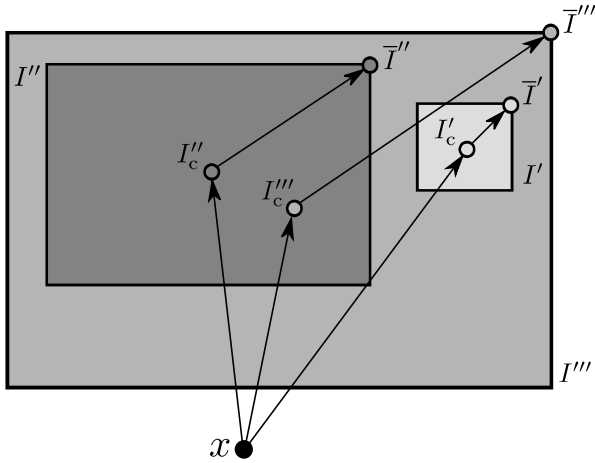


Fig. 7: The choice of  $\rho$  affects the outcome of the metric when comparing different intervals. For lower values of  $\rho$ , interval  $I'''$  is the nearest since its center  $I_c'''$  is the nearest to  $x$ . By gradually increasing  $\rho$ , 13 will select interval  $I''$ . Finally, for  $\rho$  approaching 1,  $I'$  will be the nearest.

In conclusion, when risk perception and trajectory prediction are suboptimal, the distance metric selects an action whose predicted arrival states differ the least from a known visited state.

2) *Evasion metric*: The evasion metric evaluates trajectories based on the current belief in the composition of the state space. Consider for example a state space partitioned in a safe region  $R_{safe}$ , a fatal region  $R_{fat}$ , and an uncertain region  $R_{unc}$ . More regions could be defined as long as they can be ordered from the safest to the most fatal. Each control action sequence generates bounding trajectories  $\tau$ , which OptiSHERPA computes as a succession of intervals in time. First, among all bounding trajectories, discard those that at any given time will entirely be comprised in a fatal region: for OptiSHERPA this is equivalent to generating one interval entirely composed of fatal states. Second, discard those trajectories that do not end up in the safe region at the end of the trajectory, unless no such trajectory exists. The remaining trajectories are those that have the most probability of both avoiding fatal occurrences and of restoring safety. Third, consider how the trajectories overlap the regions, e.g.  $R_{safe}$ ,  $R_{fat}$ ,  $R_{unc}$ . Assuming that all trajectories in the bounding trajectory are equally likely to occur, the bigger volume of the intersection with a fatal region, the higher the chance of a fatal occurrence when following the trajectory. Consider then two trajectories with similarly sized intersection with the fatal region, but one of which has a higher intersection with the uncertain region than the other. Since part of the uncertain states could actually be fatal, it is preferable to follow the second trajectory, reducing such a risk. This procedure can be summarized as follows:

- 1) Define a hierarchy of regions  $R_i$  where  $R_i$  has a higher probability of containing fatal states than  $R_j$ ,  $i < j$ ;
- 2) assign a weight  $w_i > 0$  to each region, so that if  $i < j$ ,  $w_i > w_j$ ;
- 3) remove from the set of feasible bounding trajectories  $\tau$  those  $\tau$  that intersect a fatal region entirely;

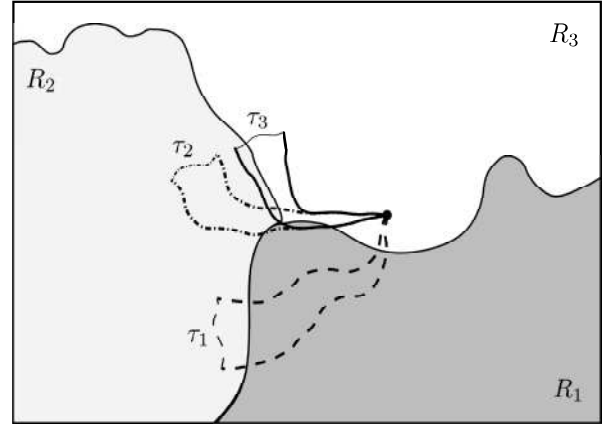


Fig. 8:  $\mathcal{S}$  comprises fatal region  $R_1$ , unknown region  $R_2$ , and safe region  $R_3$ . All  $\tau$  have a non-empty intersection with  $R_1$ . However,  $\tau_1$  crosses  $R_1$  in its entirety and is therefore the least safe;  $\tau_2$  and  $\tau_3$  also cross the region, but  $\tau_3$  has the minimal overlap with both  $R_1$  and  $R_2$ , and the biggest overlap with  $R_3$ : it is therefore the safest.

- 4) remove from  $\tau$  those  $\tau$  that do not end up in the SSS, unless this depletes  $\tau$ ;
- 5) for each  $\tau_j \in \tau$  compute the volume of the intersection  $\rho_{ij}$  with region  $R_i$ ;
- 6) the optimal control sequence  $u^*$  is the one generating  $\tau^* = \operatorname{argmin}_{\tau_j \in \tau} \sum_{R_i} \rho_{ij} \cdot w_i$ ;
- 7) apply the first element of  $u^*$ ;
- 8) if  $W = 1$ , go to 1; otherwise exit.

Fig. 8 provides an example of the above procedure. The result of applying the evasion metric is that the agent will reach for a safe region where danger is no longer perceived, while minimizing the probability of encountering a fatal occurrence.

### C. OptiSHERPA

Fig. 9 shows schematically the implementation of OptiSHERPA. At the start of the exploration, the agent is in state  $x_0$ . It will be assumed that  $W(x_0) = 0$ . OptiSHERPA generates an array of control sequences, including  $u^\pi$  given by the policy. For each sequence, the algorithm checks if the corresponding trajectory  $\tau$  is included in the SSS. Those  $\tau$  for which this is not true are removed from the array; the rest are evaluated with the distance metric. The control  $u$  corresponding to the optimal  $\tau$  is selected, and the first action of  $u^*$  is then applied. In the new state  $x_1$ , if  $W(x_1) = 0$ , the SSS is augmented with all the states currently in reach of the risk perception, and the process is repeated. If  $W(x_1) \neq 0$ , the SSS is not modified, and the evasion metric is implemented. OptiSHERPA generates an array of control sequences and corresponding  $\tau$  as before. According to the procedure of IV.B.2, an optimal control sequence  $u^*$  is found, and its first action is implemented. The evasion metric is used until  $W(x(t)) = 0$ , after which the distance metric is reimplemented. This procedure is followed until the end of exploration.

### D. Elevator control task

1) *Model*: The goal of the application is to simulate on-line training for a RL agent on-board of a fighter aircraft, which is

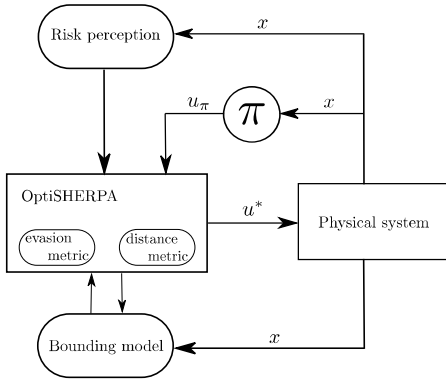


Fig. 9: The architecture of OptiSHERPA. At each time-step, the algorithm receives information in the form of current risk perception and suggested policy control  $u_\pi$ . The bounding model is used to predicted trajectories, which in turn are evaluated by implementing either the distance or the evasion metric depending on current risk perception. Then the first action of  $u^*$  is taken, which generates a new state.

constrained to fly in a range of  $\pm 40$  ft from starting height. It is assumed that a stall occurs if the angle of attack  $\alpha < -15^\circ$  or  $\alpha > 12^\circ$ . Violating this envelope causes a fatal occurrence in this environment. A bounding model of the aircraft is available as in (15) and (16):

$$\dot{h} = V \cdot \sin(\theta - \alpha) \quad (15)$$

$$\begin{pmatrix} \dot{\theta} \\ \dot{\alpha} \\ \dot{q} \end{pmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & c_{\alpha\alpha} & c_{\alpha q} \\ 0 & c_{q\alpha} & c_{qq} \end{bmatrix} \cdot \begin{pmatrix} \theta \\ \alpha \\ q \end{pmatrix} + \begin{bmatrix} 0 \\ c_{\alpha\delta_e} \\ c_{q\delta_e} \end{bmatrix} \cdot (\delta_e + \Delta\delta_e) + \begin{pmatrix} 0 \\ 0 \\ \Delta\dot{q} \end{pmatrix} \quad (16)$$

with

$$\begin{aligned} V &\in [450 \ 550] \frac{\text{ft}}{\text{s}}; c_{\alpha\alpha} \in [-0.70 \ -0.58]; \\ c_{\alpha q} &\in [0.76 \ 0.95]; c_{q\alpha} \in [-1.72 \ -1.41]; \\ c_{qq} &\in [-0.97 \ -0.79]; \Delta\delta_e \in [-0.15^\circ \ 0.15^\circ]; \\ \Delta\dot{q} &\in [-0.5 \ 0.5] c_{\alpha\delta_e} = -1.4 \cdot 10^{-3}; c_{q\delta_e} = -0.1137 \end{aligned}$$

where  $V [ft/s]$  is the constant flight speed, height  $h [ft]$  is the change in height,  $\theta [rad]$ ,  $\alpha [rad]$  are changes of pitch angle and angle of attack, and  $q [rad/s]$  is the pitch rate. Vector  $(h, \theta, \alpha, q)^T$  is the state of the system and  $\delta_e [deg]$  is the elevator deflection.

The intervals of (15) and (16) are obtained as follows. First, the nonlinear dynamics are linearised at 15 000 ft of height and  $500 \frac{\text{ft}}{\text{s}}$  of speed. Then, resulting coefficients  $c_{\alpha\alpha}$ ,  $c_{q\alpha}$  and  $c_{qq}$ , as well as speed  $V$ , are altered by  $\pm 10\%$ , whereas  $c_{\alpha q}$  is reduced between 80% and 100% of the original value.  $\Delta\dot{q}$  represents uncertainty in pitch dynamics, and  $\Delta\delta_e$  represents an error in effective deflection. Finally, terms  $c_{\alpha\delta_e}$  and  $c_{q\delta_e}$  are crisp elevator coefficients. Among all possible representations of the dynamics, the one with coefficients:

$$\begin{aligned} V &= 550 \frac{\text{ft}}{\text{s}}; c_{\alpha\alpha} = -0.58; c_{\alpha q} = 0.83; \\ c_{q\alpha} &= -1.586; c_{qq} = -0.97; \Delta\dot{q} = -0.5 \end{aligned}$$

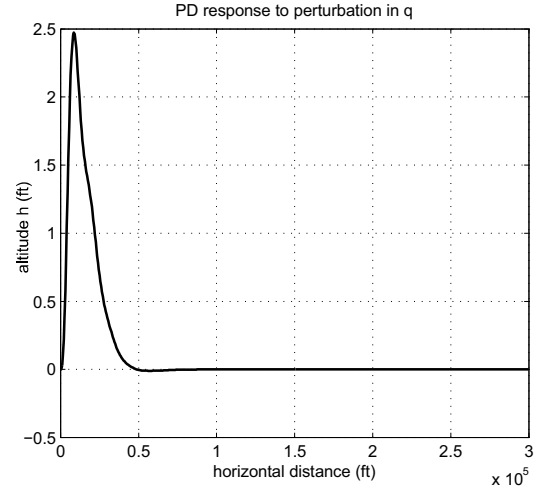


Fig. 10: The change in flight height of the system equipped with the PD controller after a perturbation in  $q$ . The final error in height is brought to zero.

is selected: cross terms  $c_{q\alpha}$  and  $c_{\alpha q}$  as the mean of their interval, and remaining terms as one of their extrema. The deflection error  $\Delta\delta_e$  is treated as noise.

A PD controller provides a baseline policy for the RL agent. Deflection  $\delta_{PD}$  is the sum of a proportional and derivative term in  $\gamma$ , of a proportional term in height  $h$ , and a proportional damping term in pitch rate:

$$\begin{aligned} \bar{\delta}_{PD} &= \frac{1}{200} \cdot h + \frac{1}{100} \cdot \dot{h} + 10 \cdot (\theta - \alpha) + 4 \cdot (\dot{\theta} - \dot{\alpha}) + \frac{1}{2} \cdot q \\ \delta_{PD} &= \max(\min(\bar{\delta}_{PD}, 1^\circ), -1^\circ) \end{aligned}$$

The controller is hand tuned to achieve a satisfactory performance in nominal conditions. Fig. 10 shows a stable and well damped response to a perturbation in  $q$  of  $0.5^\circ/s$ . The proportional term in  $h$  brings the final error in altitude to zero.

The RL agent follows an  $\epsilon$ -greedy policy by selecting random actions  $\delta \in [-1; 1]$  with probability  $\epsilon = 0.2$ , and  $\delta = \delta_{PD}$  otherwise.

With the system differing from nominal conditions, and with the addition of random deflections, the agent is not always able to abide to the constraints in height and angle of attack. A dedicated SHERPA could facilitate the task. However, with respect to the simple quadrotor model of the previous section, the aircraft model represents an increased challenge for SHERPA. In particular, consider the height dynamics represented by (15). Change in height is achieved via the flight path angle  $\gamma = \theta - \alpha$ , and in turn  $\dot{\gamma} = \dot{\theta} - \dot{\alpha} \cong 0.062 \cdot q + 0.56 \cdot \alpha$ , so that  $\dot{\gamma}$  and  $\dot{h}$  are influenced by the angle of attack. This represents the direct correlation between lift and flight path angle:  $\alpha$  must increase (decrease) in order to increase (decrease)  $\gamma$ . The dynamic in  $\gamma$  is considerably slower than the one in  $\theta$ ,  $\alpha$  and  $q$ , and the evolution in  $h$  for two different control sequences can be appreciated only after a certain amount of time-steps. Meaningful variations in the flight height can be observed only in further-time estimates, which are also the most uncertain. This motivates the use of OptiSHERPA.

The algorithm is initialized with the knowledge of the two risks in  $h$  and  $\alpha$ , and two corresponding risk perception ranges

of 15ft and  $5^\circ$  respectively. Whenever OptiSHERPA finds a range of height or angle of attack to be safe, it immediately adds this range to its internal SSS. The SSS is then an interval bounded in  $h$  and  $\alpha$  by the highest and lowest values obtained from risk-perception, and unbounded in  $\theta$  and  $q$ , which have no risk.

The agent generates one command sequence: either a full PD predicted sequence, or a random action plus the command predicted by the PD at later steps. Each sequence has a duration of 6 time-steps of 0.2s. Then, if danger is perceived, or if  $|\gamma| > 0.5^\circ$ , the agent enables OptiSHERPA. When enabled, the algorithm generates multiple command sequences, each lasting 6 time-steps, in the form  $\mathbf{u} = (u_0, u_b, \dots, u_b)$ ,  $u_0, u_b \in \{-0.75^\circ, -0.25^\circ, 0^\circ, 0.25^\circ, 0.75^\circ\}$ . Command  $u_0$  is the OptiSHERPA proposed action, and the constant  $u_b$  acts as a backup. This formulation of  $\mathbf{u}$  reduces the total number of metric evaluations from  $5^6$  to 25. The distance (14) is computed for all the above sequences, with  $\mathbf{v}_r = [0.2(\text{ft})^{-1}, \frac{36}{\pi}, \frac{36}{\pi}, \frac{4}{\pi}\text{s}]^T$ ,  $\rho = 0.5$ ,  $\lambda = 0.1$ .

If danger is perceived, the evasion metric is implemented, and the state space is partitioned into distinct regions, depending on the risk perception and the knowledge of the SSS. If only one of either  $\text{RFSS}_h$  or  $\text{RFSS}_\alpha$  is detected, the following are defined:

- a “black” region  $R_{\text{bk}}$  comprising all states with fatal values of  $h$  (or  $\alpha$ );
- a “white” region  $R_{\text{wh}}$  comprising all states within risk perception of the nearest fatal value;
- a “green” region  $R_{\text{gn}}$  equal to the  $\text{SSS} \setminus R_{\text{wh}}$ ;
- a “grey” region  $R_{\text{gy}}$  equal to the remainder of  $\mathcal{S}$ .

The white region is introduced in order to encourage the metric to move away from those regions where risk is still perceived, thus exiting the evasion cycle. If both  $h$  and  $\alpha$  are perceived, then two regions are added: a “red” region  $R_r$  comprising the states that are fatal for one feature, but not for the other, and a “blue region”  $R_{\text{bl}}$  comprising those states within risk perception of the nearest fatal value of one feature, but not of the other. Finally, to each region is assigned a weight

$$w_{\text{bk}} = 2, w_r = 0, w_{\text{gy}} = -0.1, w_w = -0.2, w_{\text{bl}} = -1, w_{\text{gn}} = -2.$$

A total of 500 trials is performed from the starting perturbed condition  $h = \alpha = \theta = 0$ ,  $q = 0.5^\circ$ . Each trial has a duration of 600 timesteps, equivalent to 120s. In each trial, both the “original” RL agent and the one augmented with SHERPA are run at the same time. The added noise on the elevator deflection and the random action were the same at all time for both the original and the augmented agent, albeit SHERPA is allowed to dismiss the random action as previously exposed.

Fig. 11 shows a typical result for the task. The upper plot shows the two trajectories with the application of OptiSHERPA and with the original agent. At the start of the task, the original agent manages to keep deviations of  $\gamma$ ,  $\alpha$ , and  $q$  within reasonable limits, thus satisfying the requirements for OptiSHERPA. However, the effect of OptiSHERPA becomes noticeable when a risk in  $h$  is perceived. In all those instances, the algorithm finds a reactive solution in the first iterations, increasing the pitch rate and the angle of attack. Then, in later iterations when recovery is incipient, OptiSHERPA applies a

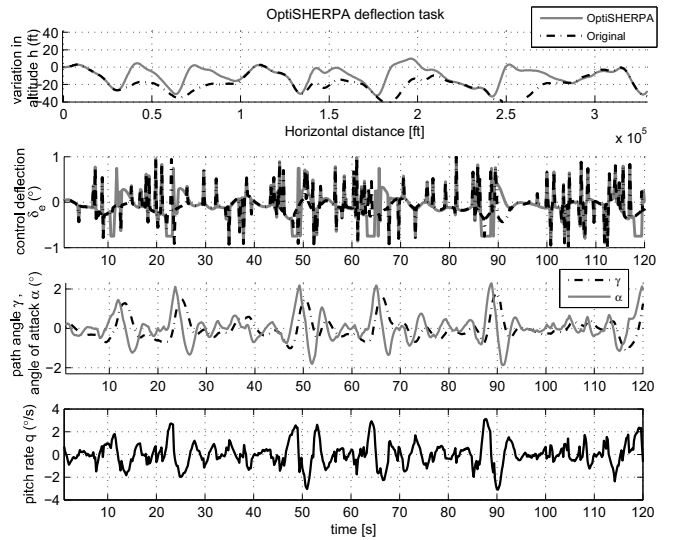


Fig. 11: The top plot represents the trajectory with the original RL agent and with the addition of OptiSHERPA. The 2<sup>nd</sup> plot represents deflection  $\delta_e$  for the two different agents; OptiSHERPA selects constant maximal deflections during recovery by optimising the evasion metric. The 3<sup>rd</sup> shows the variations of  $\gamma$  and  $\alpha$  during flight with OptiSHERPA. The bottom plot shows pitch rate  $q$ .

deflection of opposite sign. The reason for this is as follows. Both the flight path angle and the angle of attack have changed sign, so that the bounding model predicts a recovery in flight height in near time. However, the positive pitch angle could increase further the angle of attack. Since such levels have not been experimented before, the evasion metric considers this a potentially threatening situation. Thus, the best course of action for OptiSHERPA is to achieve negative pitch rate to reduce  $\alpha$ . Therefore, lift increases during the first iterations, and reduces it after recovery is in progress. As for  $q$ , rapid changes are performed during recoveries; when not in recovery, the agent keeps  $q$  within acceptable limits. By means of following the evasion metric, OptiSHERPA consistently recovers in a more incisive and fast way than the equivalent PD action, and prevents violations in height during the whole flight. Without this additional recovery capability, the original RL agent violates the constraint in flight height in 53% of the trials. It is necessary to consider, however, that the original controller is always adopting an exploratory  $\epsilon$ -greedy policy. A direct comparison between the two should therefore be avoided.

The angle of attack plays a crucial role in determining whether a recovery can be achieved with the limited distance in feet allowed by the risk perception. Fig. 12 shows one such case. Prior to the violation, a series of positive random deflections reduce considerably the flight path angle, while at the same time decreasing the angle of attack and generating negative pitch rate. As soon as the risk is perceived, and OptiSHERPA is engaged with the evasion metric, the agent increases the pitch rate and starts generating additional lift through  $\alpha$ , but the violation still occurs. Nonetheless, it can be seen how OptiSHERPA suggest the correct maximum

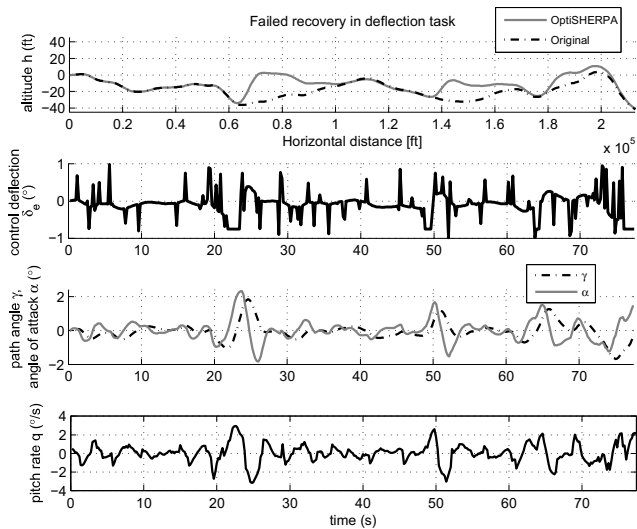


Fig. 12: A typical violation of the flight height constraints.

deflection, and how the residual negative  $\gamma$  is reduced in the instants prior to violation. Approximately 10% of the simulations incurred in such a violation for the agent with OptiSHERPA. However, two considerations must be added. First, all violations occurred following cumulative random deflections, giving rise to conditions from which recovery was very difficult even with maximum deflection, and resulted in minor violations (as in Fig. 12). Second, all cases where violation occurred for OptiSHERPA were either concurrent or preceded by a violation for the original RL agent.

OptiSHERPA showed the correct behavior in case of imminent flight height violation, selecting strong deflections in the first iterations and subsequently reducing the variation of angle of attack introduced. This resulted in quick and safe recoveries. For a fraction of the runs, the stochasticity in the agent action gave rise to small flight height violations, although without any reduction in performance when compared to the alternative original RL agent, which was equally affected.

## V. CONCLUSIONS

This work presents a new approach for autonomous agents in dangerous environments. The presence of fatal and lead-to-fatal states constitutes the motivation for the approach, with risk perception as its main assumption. The notions of safe control and backup are introduced to present an algorithm for safe exploration: SHERPA. By relying on a bounding model of the dynamics, SHERPA allows those policy actions for which the system can be brought near a previously known state by satisfying a closeness condition, thus promoting safety. SHERPA is tested in one simulated quadrotor application, achieving safety. Subsequently, OptiSHERPA is introduced in order to handle tasks for which either the risk perception or the bounding model were insufficient for the original SHERPA. These limitations are addressed by adding metrics, which provide the agent with informed options, and by explicitly including an evasion strategy in those cases where danger is

imminent. OptiSHERPA is tested on a second simulated task: maintaining straight flight for a fighter aircraft exploring a strict envelope, with the addition of noise, uncertain dynamics and random exploratory actions. The application shows how the resulting RL agent with OptiSHERPA manages reliable control, adopting a very reasonable behavior during recoveries. The proposed approach constitutes a significant effort into tackling the exploration problem for RL agents on a general level, while not focusing on a particular category of tasks. This is reflected in the absence, within the algorithm, of a high-fidelity model for exploration. This strategy is therefore in line with the model-free approach of Reinforcement Learning and of adaptive controllers in general. Future development will include an investigation of: risk perception for real-life scenarios by using sensor-information, methods for the autonomous and possibly adaptive selection of the open parameters of the algorithms (such as  $m$  and  $v_r$ ) and representations of uncertainty not based on interval methods to reduce computational complexity.

## REFERENCES

- [1] H. Kwakernaak, "Robust control and  $h_\infty$ -optimization - tutorial paper," *Automatica*, vol. 29, no. 2, pp. 255–273, 1993. [Online]. Available: <http://doc.utwente.nl/29962/>
- [2] E. van Kampen, "Continuous adaptive critic flight control aided with approximated plant dynamics." Reston, Va: American Institute of Aeronautics and Astronautics, 2006 2006.
- [3] S. Ferrari and R. F. Stengel, "Online adaptive critic flight control," *J. Guid. Control Dyn.*, vol. 27, no. 5, pp. 777–786, 2004.
- [4] R. S. Sutton and A. G. Barto, *Introduction to Reinforcement Learning*, 1st ed. Cambridge, MA, USA: MIT Press, 1998.
- [5] D. Haussler, "Probably approximately correct learning," in *Proc. of the 8th Nat. Conf. Artif. Intell. - Vol. 2*, ser. AAAI'90. AAAI Press, 1990, pp. 1101–1108.
- [6] D. Zhao and Y. Zhu, "Meca near-optimal online reinforcement learning algorithm for continuous deterministic systems," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 2, pp. 346–356, Feb 2015.
- [7] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming*, 1st ed. Athena Scientific, 1996.
- [8] J. Si and Y.-T. Wang, "Online learning control by association and reinforcement," *IEEE Trans. Neural Netw.*, vol. 12, no. 2, pp. 264–276, Mar 2001.
- [9] R. Bellman, *Dynamic Programming*, 1st ed. Princeton, NJ, USA: Princeton University Press, 1957.
- [10] M. C. Choy, D. Srinivasan, and R. L. Cheu, "Neural networks for continuous online learning and control," *IEEE Trans. Neural Netw.*, vol. 17, no. 6, pp. 1511–1531, Nov 2006.
- [11] Y. J. Liu, L. Tang, S. Tong, C. L. P. Chen, and D. J. Li, "Reinforcement learning design-based adaptive tracking control with less learning parameters for nonlinear discrete-time mimo systems," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 1, pp. 165–176, Jan 2015.
- [12] Y.-J. Liu, Y. Gao, S. Tong, and Y. Li, "Fuzzy approximation-based adaptive backstepping optimal control for a class of nonlinear discrete-time systems with dead-zone," *IEEE Trans. Fuzzy Syst.*, vol. 24, no. 1, pp. 16–28, 2016.
- [13] D. Liu, X. Yang, D. Wang, and Q. Wei, "Reinforcement-learning-based robust controller design for continuous-time uncertain nonlinear systems subject to input constraints," *IEEE Trans. Cybern.*, vol. 45, no. 7, pp. 1372–1385, 2015.
- [14] A. Arleo, F. Smeraldi, and W. Gerstner, "Cognitive navigation based on nonuniform gabor space sampling, unsupervised growing networks, and reinforcement learning," *IEEE Trans. Neural Netw.*, vol. 15, no. 3, pp. 639–652, May 2004.
- [15] C. Wang, Y. Li, S. S. Ge, and T. H. Lee, "Optimal critic learning for robot control in time-varying environments," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 10, pp. 2301–2310, Oct 2015.
- [16] T. Shimizu, R. Saegusa, S. Ikemoto, H. Ishiguro, and G. Metta, "Robust sensorimotor representation to physical interaction changes in humanoid motion learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 5, pp. 1035–1047, May 2015.

- [17] M. Pecka and T. Svoboda, *Modelling and Simulation for Autonomous Systems: First International Workshop, MESAS 2014, Rome, Italy, May 5-6, 2014, Revised Selected Papers*. Cham: Springer International Publishing, 2014, ch. Safe Exploration Techniques for Reinforcement Learning – An Overview, pp. 357–375.
- [18] M. Heger, “Consideration of risk in reinforcement learning,” in *Proc. of the 11th Int. Conf. Mach. Learn.*, 1994, pp. 105–111.
- [19] J. García and F. Fernández, “A comprehensive survey on safe reinforcement learning,” *J. Mach. Learn. Res.*, vol. 16, pp. 1437–1480, 2015.
- [20] P. Thomas, G. Theodorou, and M. Ghavamzadeh, “High confidence policy improvement,” in *Proc. of the 32nd Int. Conf. Mach. Learn. (ICML-15)*, D. Blei and F. Bach, Eds. JMLR Workshop and Conference Proceedings, 2015, pp. 2380–2388.
- [21] R. Neuneier and O. Mihatsch, “Risk sensitive reinforcement learning,” in *Advances in Neural Information Processing Systems*, 1999, pp. 1031–1037.
- [22] O. Mihatsch and R. Neuneier, “Risk-sensitive reinforcement learning,” *Machine Learning*, vol. 49, no. 2, pp. 267–290.
- [23] Y. Shen, M. J. Tobia, T. Sommer, and K. Obermayer, “Risk-sensitive reinforcement learning,” *Neural Computation*, vol. 26, no. 7, pp. 1298–1328, 2014.
- [24] P. Geibel and F. Wysotzki, “Risk-sensitive reinforcement learning applied to control under constraints,” *J. Artif. Intell. Res.*, vol. 24, pp. 81–108, 2005.
- [25] S. P. Coraluppi and S. I. Marcus, “Risk-sensitive and minimax control of discrete-time, finite-state markov decision processes,” *Automatica*, vol. 35, no. 2, pp. 301 – 309, 1999. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0005109898001538>
- [26] K. Driessens and S. Džeroski, “Integrating guidance into relational reinforcement learning,” *Machine Learning*, vol. 57, no. 3, pp. 271–304, 2004.
- [27] B. D. Argall, S. Chernova, M. Veloso, and B. Browning, “A survey of robot learning from demonstration,” *Robotics and Autonomous Systems*, vol. 57, no. 5, pp. 469–483, 2009.
- [28] D. Martinez, G. Alenya, and C. Torras, “Safe robot execution in model-based reinforcement learning,” in *2015 IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Sept 2015, pp. 6422–6427.
- [29] C. Gehring and D. Precup, “Smart exploration in reinforcement learning using absolute temporal difference errors,” in *Proc. 2013 Int. Conf. Autonomous Agents Multi-agent Syst.* International Foundation for Autonomous Agents and Multiagent Systems, 2013, pp. 1037–1044.
- [30] A. Hans, D. Schneegaß, A. M. Schäfer, and S. Udluft, “Safe exploration for reinforcement learning,” in *ESANN*, 2008, pp. 143–148.
- [31] M. Pecka, K. Zimmermann, and T. Svoboda, “Safe exploration for reinforcement learning in real unstructured environments.”
- [32] F. J. G. Polo and F. F. Rebollo, “Safe reinforcement learning in high-risk tasks through policy improvement,” in *Adaptive Dynamic Programming And Reinforcement Learning (ADPRL), 2011 IEEE Symposium on*, April 2011, pp. 76–83.
- [33] T. M. Moldovan and P. Abbeel, “Safe Exploration in Markov Decision Processes,” in *Proc. 29th Int. Conf. Mach. Learn.* icml.cc / Omnipress, 2012, p. 188.
- [34] L. P. Kaelbling, M. L. Littman, and A. W. Moore, “Reinforcement learning: A survey,” *J. Artif. Intell. Res.*, vol. 4, pp. 237–285, 1996.
- [35] C. Szepesvári, “Algorithms for reinforcement learning,” *Synthesis Lectures Artif. Intell. Mach. Learn.*, vol. 4, no. 1, pp. 1–103, 2010.
- [36] P. Geibel, “Reinforcement learning with bounded risk,” in *ICML*, 2001, pp. 162–169.
- [37] T. Fraichard and H. Asama, “Inevitable collision states a step towards safer robots?” *Advanced Robotics*, vol. 18, no. 10, pp. 1001–1024, 2004.
- [38] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, no. 3, pp. 279–292, 1992.
- [39] R. E. Moore, *Interval analysis*. Prentice-Hall Englewood Cliffs, 1966, vol. 4.
- [40] O. Khatib, “Real-time obstacle avoidance for manipulators and mobile robots,” *Int. J. Robot. Res.*, vol. 5, no. 1, pp. 90–98, 1986.
- [41] K. Ok, S. Ansari, B. Gallagher, W. Sica, F. Dellaert, and M. Stilman, “Path planning with uncertainty: Voronoi uncertainty fields,” in *2013 IEEE Int. Conf. Robot. Autom. (ICRA)*. IEEE, 2013, pp. 4596–4601.
- [42] S. G. Loizou, H. G. Tanner, V. Kumar, and K. J. Kyriakopoulos, “Closed loop motion planning and control for mobile robots in uncertain environments,” in *Proc. 42nd IEEE Conf. Decision Control, 2003.*, vol. 3. IEEE, 2003, pp. 2926–2931.
- [43] V. J. Lumelsky and A. A. Stepanov, “Path-planning strategies for a point mobile automaton moving amidst unknown obstacles of arbitrary shape,” *Algorithmica*, vol. 2, no. 1-4, pp. 403–430, 1987.
- [44] K. P. Tee, S. S. Ge, and E. H. Tay, “Barrier lyapunov functions for the control of output-constrained nonlinear systems,” *Automatica*, vol. 45, no. 4, pp. 918–927, 2009.
- [45] P. Wieland and F. Allgöwer, “Constructive safety using control barrier functions,” *IFAC Proceedings Volumes*, vol. 40, no. 12, pp. 462–467, 2007.
- [46] M. Z. Romdlony and B. Jayawardhana, “Uniting control lyapunov and control barrier functions,” in *53rd IEEE Conf. Decision Control*. IEEE, 2014, pp. 2293–2298.
- [47] T. M. Lam, H. W. Boschloo, M. Mulder, and M. M. Van Paassen, “Artificial force field for haptic feedback in uav teleoperation,” *IEEE Trans. Syst., Man, Cybern.A, Syst., Humans*, vol. 39, no. 6, pp. 1316–1330, 2009.



**Tommaso Mannucci** obtained his BSc-degree in Aerospace Engineering at University of Pisa in 2009 and his MSc-degree in Aerospace Engineering and Flight Dynamics in 2012 with a thesis on numerical optimization for attitude recovery. Since 2013 he is a PhD candidate at the Control and Simulation division of the Aerospace Faculty at Delft University of Technology. His main research interests are Adaptive Control, Machine Learning and Reinforcement Learning for UAVs.



**Dr. Erik-Jan van Kampen** obtained his BSc-degree in Aerospace Engineering at Delft University of Technology in 2004, his MSc-degree in Control and Simulation in 2006, and a PhD-degree in Aerospace Engineering in 2010. He is currently an assistant professor at the Control and Simulation division at Delft University of Technology. His research interests are Intelligent Flight Control, adaptive control, and interval optimization.



**Dr. Cornelis de Visser** received the M.Sc. degree from the Delft University of Technology in 2007. In 2011 he received his Ph.D degree from the faculty of Aerospace Engineering at the Delft University of Technology in The Netherlands. Currently, he is an assistant professor at the Control and Simulation section of the Delft University of Technology. His research interests include aircraft system identification, flight envelope prediction, fault tolerant control, and multivariate spline theory.



**Dr. Q.P. (Qiping) Chu** received his Ph.D. degree from the Faculty of Aerospace Engineering, Delft University of Technology, The Netherlands, in 1987. Currently, he is an Associate Professor at the Faculty of Aerospace Engineering, Delft University of Technology and the head of aerospace guidance, navigation and control cluster within the section of control and simulation. He has (co)authored more than 200 journal and conference papers ranged from adaptive control, nonlinear control, robust control and intelligent control to nonlinear state estimation, system identification and nonlinear optimization for aerospace vehicles.