

# Intelligent Behavior Utilizing Light-weight Diagnosis

Martin Zimmermann<sup>1</sup>, Franz Wotawa<sup>1</sup>, and Ingo Pill<sup>1</sup>

<sup>1</sup>Affiliation not available

November 1, 2021

## Abstract

Intelligence in its decisions is a trait that we have grown to expect from a cyber-physical system. In particular that it makes the right choices at runtime, i.e., those that allow it fulfill its tasks, even in case of faults or unexpected interactions with its environment. Analyzing how to continuously achieve the currently desired (and possibly continuously changing) goals and adapting its behavior to reach these goals is undoubtedly a serious challenge. This becomes even more challenging if the atomic actions a system can implement become unreliable due to faulty components or some exogenous event out of its control. In this paper, we propose a solution for the presented challenge. In particular, we show how to adopt a light-weight diagnosis concept to cope with such situations. The approach is based on rules coupled with means for rule selection that are based on previous information regarding the success or failure of rule executions. We furthermore present a Java-based framework of the light-weight diagnosis concept, and discuss the results obtained from an experimental evaluation considering several application scenarios. At the end, we present a qualitative comparison with other related approaches that should help the reader decide which approach works best for them.

Corresponding author(s) Email: [wotawa@ist.tugraz.at](mailto:wotawa@ist.tugraz.at)

## ToC Figure

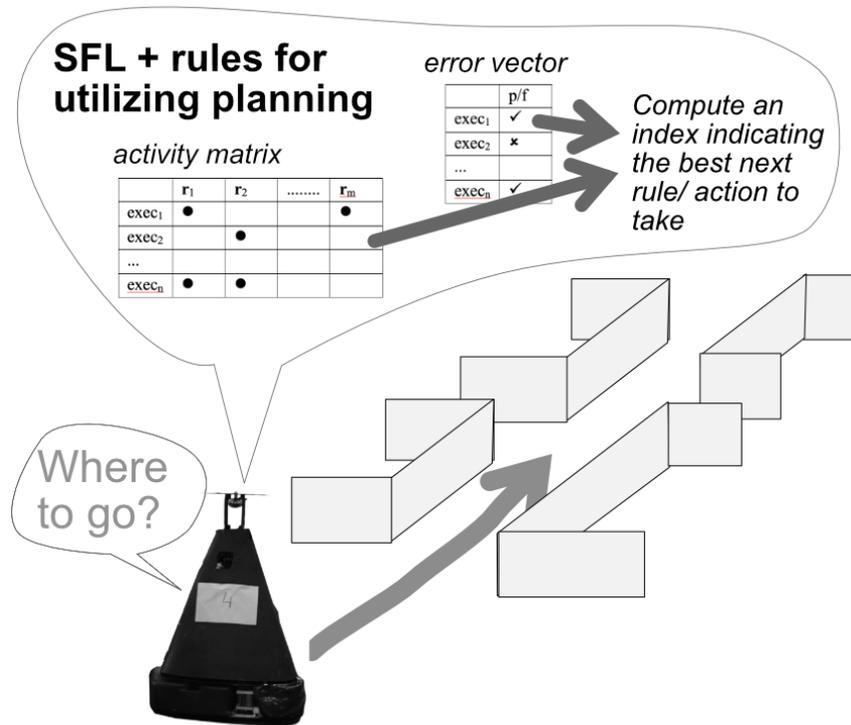


Figure 1: **ToC Figure.** This article deals with utilizing SFL in combination with rule-based action execution for controlling smart agents in an uncertain environment. The focus is on providing an intelligent agent with efficient capabilities for fulfilling given goals in case of missing information or faults. In addition, to foundations the article comprises an experimental evaluation, and a qualitative comparison to other related methods.

## Introduction

Be it our smartphones having to last through the day, automated industry plants, or autonomous robots and cars: We expect all these systems to make intelligent decisions to effectively and efficiently perform their tasks, regardless of encountered issues and changes in the environment.

There is a variety of techniques that we can draw on for implementing such behavior, including calculi like the situation calculus (Boutilier et al., 2000; Gspandl et al., 2011) or dynamic planning concepts that allow us to react to changes in the environment (Connell and La, 2017). Regardless of the technique, in principle we are searching for action sequences that allow us to achieve our current goals. In this context, we certainly not only have to check whether we correctly perceive and assess the environment, but in practice, we're also likely to suffer from unreliable actions. With our work as presented in this paper, we are focusing on the latter.

Reasoning about an issue's exact origin(s), e.g., with model-based diagnosis (MBD) (Reiter, 1987; de Kleer and Williams, 1987), allows us to search for an ideal mitigation strategy. In practice, however, (a) there is seldom enough data to precisely isolate a problem's source(s) so that we end up with a set of candidates, and

(b) the reasoning’s complexity might prohibit us from making fast decisions. Moreover, in MBD we require in addition to observations of the system a system model that captures the behavior sufficiently to allow deriving diagnosis candidates.

Complementing detailed and complete diagnosis concepts like MBD, in recent years spectrum-based fault localization (SFL) (Abreu et al., 2009b) has been gaining in attention. With SFL we evaluate *execution data* about which component was involved in this or that observed behavior. The result is a ranking how suspicious the individual components are to have caused the failing behaviors. Traditionally, SFL has been employed for software debugging, but, e.g., it was shown recently in (Pill and Wotawa, 2018) how to translate the idea for a *static* diagnosis of knowledge-bases used in automated reasoning. As will be discussed in the next section, there we observed which of the knowledge-base’s rules were involved in the individual reasoning tasks and whether the tasks were successful. Because of the ease of implementing SFL and the smaller number of information required for diagnosis, SFL presents a representative of light-weighted diagnosis methodologies.

We elaborate on such related work in that we translate and extend the basic SFL concept to accommodate also *dynamic*, live settings. Our concept thus allows us to continuously evaluate the success-rates of system actions via considering the success/failing of previously executed action *sequences* a.k.a. plans. We continuously update a corresponding reliability measure for *each* action, and use these data to (a) re-evaluate currently implemented action plans, and (b), when reasoning about new plans for achieving future goals. We require only very limited data for our concept, i.e., which *sequences* failed or succeeded in the past and which components were involved. From a technical perspective, we describe a system’s actions via specific rules, pre- and postconditions and use these data for our reasoning.

While we do not isolate an issue’s exact source(s) (like with MBD), we will show in our evaluation that our compromise between preciseness and computational complexity allows us to *dynamically, effectively* and *efficiently* cope with faults and other events that result in unreliable actions. As we will discuss, we rely solely on SFL for our reasoning and thus passively observed executions without deploying any exploratory component like in reinforcement learning.

## Preliminaries

MBD is undoubtedly a powerful technique for precisely isolating a problem’s origin(s). We need a special model for this reasoning though, and while MBD is complete with respect to this model, the entailed computations can become quite complex. That is, the diagnosis search space is exponential in the number of the health state variables that we have to introduce. The more health state variables we have, the more faults can be found, but the larger the search space is. Usually, we end up with more than one diagnosis matching the data, so that we have to choose one as a working hypothesis.

With SFL, we take a different approach and consider the involvement of components in failing *and* passing behavior. The reasoning then follows the idea that some component that is always involved in faulty behavior but never in correct behavior is very suspicious of being the source of the troubles (and vice versa). Since components are usually involved in both faulty *and* correct behavior, as well as the possibility that some faulty components cancel each other out, leading to correct behavior, many *similarity coefficients*, e.g. (Jones and Harrold, 2005; Abreu et al., 2009a), for computing a component’s suspiciousness have been proposed.

An intrinsic advantage of considering multiple executions in SFL by default is that fault masking (when multiple faults lead to correct output observations) has less effects on the reasoning, that is, if the set of observed behaviors is representative enough to contain also behavior without the masking effect. If the set is indeed representative and the faults *always* mask each other, then we are possibly facing an equivalent mutant so that we might want to consider the “faults” as implementation alternatives. For their computation, we consider the corresponding execution data about which component was involved in which behavior (stored in a matrix also referred to as spectrum), and whether some behavior is violating or complying with our

expectations (the so-called error vector). Based on the components' suspiciousness values, we establish a ranking.

**Definition 1** An activity matrix or spectrum  $A$  is an  $n \times m$  matrix, where for each of the  $n$  system components we have  $m$  rows for  $m$  considered behaviors  $b_j$ . A cell  $a_{ij}$  is labeled with 1 iff component  $c_i$  is involved in  $b_j$ , otherwise with 0.

**Definition 2** An error vector  $e$  for some spectrum  $A$  (Def. 1) is a vector of length  $m$  (a  $1 \times m$  matrix) such that  $e_j = 1$  iff  $b_j$  in  $A$  violates the expectations, and  $e_j = 0$  otherwise.

From  $A$  and  $e$ , we derive for each  $c_i$  four *frequencies*  $n_{CN}(c_i)$ ,  $n_{CE}(c_i)$ ,  $n_{VN}(c_i)$  and  $n_{VE}(c_i)$  that capture in how many **C**orrect and **V**iolating behaviors (the rows) in  $A$  some  $c_i$  was **E**xecuted or **N**ot. In Tab. 1 we summarized the calculation of the four frequencies.

From these, we can compute several similarity coefficients like Ochiai, Tarantula, or Jaccard for estimating a component  $c_i$ 's suspiciousness  $D(c_i)$ :

$$\text{Ochiai} : D(c_i) = \frac{n_{VE}(c_i)}{\sqrt{(n_{VE}(c_i)+n_{VN}(c_i)) \cdot (n_{VE}(c_i)+n_{CE}(c_i))}}$$

$$\text{Tarantula} : D(c_i) = \frac{\frac{n_{VE}(c_i)}{n_{VE}(c_i)+n_{VN}(c_i)}}{\frac{n_{VE}(c_i)}{n_{VE}(c_i)+n_{VN}(c_i)} + \frac{n_{CE}(c_i)}{n_{CE}(c_i)+n_{CN}(c_i)}}$$

$$\text{Jaccard} : D(c_i) = \frac{n_{VE}(c_i)}{n_{VE}(c_i)+n_{VN}(c_i)+n_{CE}(c_i)}$$

For software, it is very easy to come up with these data. So while SFL was originally developed for that domain, it has been employed also in other contexts. In (Pill and Wotawa, 2018), it was shown, e.g., that we can exploit SFL in the context of logic reasoning with knowledge-bases. While a knowledge base is not a program that we execute in the traditional sense, one can record the *rules that are used when reasoning about a problem* and use these data to define the spectrum. The reasoning processes for individual problems with the same knowledge-base then define the  $b_j$  for  $A$ . For defining the *error vector*, it was suggested to inspect whether one would derive a contradiction *and* whether one would fail to derive the expected conclusions.

In the next section, we will show how to extend this concept to a live setting and a continuous assessment of a system's actions' *reliability*. Our aim will not be to establish a ranking about which rules fail in practice, but rather to express our confidence in the individual rules working out as expected.

There is quite a variety of tools and reasoning engine techniques that we could have adapted for evaluating our reasoning concept in practice.

Since our motivation has been to identify and reason about reliable action sequences, we focused on an available framework that encodes actions into a knowledge-base of rules and where we then reason with these rules describing the actions. This framework called RBL (for rule-based language) is available for Java programs and was proposed in (Zimmermann and Wotawa, 2020). There is also an implementation available at <https://github.com/martinzimmermann/RBL-Framework/releases/tag/CPS-RTSA>. Some interesting aspects for us were that RBL not only allows to execute sequences, but there is also some functionality to continuously (re-)design the sequence during execution.

Planning with rules is of course not a new concept introduced by RBL, but it has been studied before, e.g., in (Fikes and Nilsson, 1971; Blum and Furst, 1997; Kautz and Selman, 1996). The reason we chose

Table 1: Four *frequencies* catching how often some component  $c_i$  was involved in specific behavior:  $n_{pq}(c_i)$

$n_{CN}(c_i)$	... # of correct behaviors (C) s.t. $c_i$ was not executed (N).
$n_{CE}(c_i)$	... # of correct behaviors (C) s.t. $c_i$ was executed (E).
$n_{VN}(c_i)$	... # of violating behaviors (V) s.t. $c_i$ was not executed (N).
$n_{VE}(c_i)$	... # of violating behaviors (V) s.t. $c_i$ was executed (E).

RBL for demonstrating our concept is that already its original runtime engine allows to derive *and* executes some action sequence and it was also designed to exploit feedback from the execution in some continuous re-planning concept. This made it an ideal candidate for adopting our concept of implementing an engine that allows to make intelligent decisions where we continuously assess the situation, derive diagnostic data via a special SFL concept and derive new plans that are most promising on achieving the desired goals.

In RBL, a system’s environment is modeled by a dynamic list of corresponding beliefs whereas the system is described via *rules*. Such a rule comprises its preconditions, its postconditions, an action (Def. 4), a repair routine (Def. 5), and a weight catching our confidence in this rule’s success.

**Definition 3** A *rule*  $R = (G, P, a, r, w)$  consists of a finite set  $G$  of preconditions, a finite set  $P$  of postconditions, an action  $a$ , a repair routine  $r$  and a weight  $w$ , where  $G$  and  $P$  are non-disjoint sets. Rule  $R$  is guarded by  $G$  and can only be executed if all  $g \in G$  are (currently) known as beliefs. Iff executing  $R$  (and thus  $a$ ) is successful, all  $p \in P$  are added to (or removed from) the runtime engine’s beliefs accordingly. If it fails,  $R$ ’s repair routine  $r$  is invoked.

**Definition 4** An *action*  $a$  is a function that interacts with the environment. It returns  $\top$  (true) if the interaction was successful and  $\perp$  (false) otherwise.

**Definition 5** A *repair routine*  $r$  is a function that has to be designed by the user and repairs the runtime engine’s belief such as to reach a correct and coherent state.

**Definition 6** A *finite plan*  $\pi$  is a finite sequence of rules  $R_1, \dots, R_n$ , such that the individual rules’ preconditions are met and the desired goal is reached when executing  $\pi$ .

RBL’s Plan-Execute-Update cycle perfectly fits our concept with its three stages that we can adopt also for our SFL-based reasoning concept. As we will show in the next Section, we will consider our reasoning concept in the planning phase to generate reliable action sequences and update the spectrum in the update phase.

1. **Planning:** Search for a plan  $\pi$  with the highest chance of success (lowest costs related to a plan’s actions’ weights).
2. **Execution:** Execute a plan’s rule sequence and track the rules’ successfulness. If executing a rule  $R$ ’s associated action failed, (1)  $\pi$ ’s execution is terminated and (2) the  $R$ ’s repair routine is invoked leading to a coherent state.
3. **Update:** Update the rules’ execution history and re-calculate the rules’ weights.

While we chose RBL for our implementation, we opted to keep our concept general and did not want it to be restricted to using RBL’s custom domain-specific language. Consequently, we decided to implement the Planning Domain Definition Language (PDDL) (Ghallab et al., 1998) for our front-end, to show that our approach is available and applicable to all domains compatible with PDDL (and that it likewise can be implemented also in other frameworks).

PDDL was introduced in 1998 as a domain independent planning language compatible with many algorithms. The main components of PDDL are the *domain description* and the *problem description* (Ghallab et al., 1998). In the first, we define a set of actions with preconditions and effects. These actions usually have parameters that are populated with predicates during planning. In the problem description we describe the initial state and the goal state - also using predicates. The planning algorithm’s job is then to derive a plan (as an ordered list of actions and their corresponding parameter assignments) that when executed leads from the initial state to the goal state.

# Using Diagnostic Reasoning to Compute a System’s Actions’ Reliability and Foster Intelligent Behavior

As we outlined in the previous Section, traditionally, SFL has been deployed in a static context where we consider a test suite’s execution for an a-posteriori identification of faulty components. Our aim is quite different in that we focus on a live setting where we continuously collect new execution data. Constantly analyzing these data via diagnostic reasoning, we establish a reliability measure for each of a system’s actions (they act as our “components”), recognize failed action sequences, and when re-planning (and when deriving future plans) we aim to select the most reliable actions (technically it is their rules) for achieving our goals. That is, those rules (and their sequences) that are least likely to fail.

The three-phase cycle implemented in RBL’s run-time engine as discussed in the preliminaries perfectly fits the demands for our control-concept, so that we did not have to implement it ourselves but focused on adapting the engine for our reasoning. Also, other dynamic planning environments feature similar control concepts so that our approach could be easily adopted there needing only some adaptations to accommodate our reasoning. RBL’s engine, for instance, used a quite different cost function (a specific weight model) and related plan optimization concept, so that we needed to adapt the planning algorithm to support our own SFL-related reasoning.

But in general, once we translate the static SFL idea to a dynamic context, deploying it does not require massive changes in a corresponding run-time engine like that of RBL. That is, in principle, we need to continuously

1. derive and execute a plan  $\pi$  to achieve the current goals,
2. get feedback about  $\pi$ ’s success, and
3. compute the individual rules’ frequencies and in turn their suspiciousness/reliability to be considered when making future decisions (when deriving future plans).

Please note that we use the terms actions and rules interchangeably in this manuscript, since technically we reason with rules that describe a system’s actions.

Formally and from an abstract point of view, for adopting SFL in our dynamic context, we have to add another row to  $A$  for specifying which rules were part of  $\pi$  whenever some plan  $\pi$  failed or succeeded. Furthermore, we have to enlarge the error vector  $e$  to report also whether  $\pi$  failed or not. From  $A$  and  $e$ , we can compute similarity coefficients via the formulae depicted in the previous section. In principle, also a sliding window could be used such as to only consider recent data, which might be desirable for some dynamic applications.

In practice, the computation is less complex since we can keep track of all the rules’ four frequencies’ values, and whenever a plan  $\pi$  fails or succeeds, we increase the appropriate frequencies by one and recalculate the coefficients. In our context, we do not establish a ranking with these values, but consider them as reliability measures for the corresponding rule. This value describes how likely a component is failing, so that like for standard SFL applications we have that the lower the value, the less likely a component is to cause troubles. A plan’s reliability is then computed from the reliability of its individual actions as follows:

**Definition 7** plan  $\pi$ ’s reliability  $r(\pi)$  is computed as the sum of the plan’s individual rules’ reliability.

In Algorithm 1, we explicitly illustrate the steps needed for our concept. We show the action calls and repair routine invocations associated with a rule and illustrate the required loops and decisions. Please note that function *update\_frequencies*( $E, res$ ) serves to update the frequencies and subsequently recompute the reliability values for all the individual rules as outlined above. It has two arguments: a list  $E$  containing those rules that have been executed for plan  $\pi$ , and  $res$  encoding  $\pi$ ’s success.

**Algorithm 1:** An algorithm that supervises a plan’s execution and updates the frequencies and reliability for all rules

*Input:* a valid plan  $\Pi$  and a function *update\_frequencies* that updates the frequencies and rule coefficients

*Output:*  $\top$  if the plan execution succeeded,  $\perp$  otherwise

```

execute_plan( $\pi$ , update_frequencies)
   $E \leftarrow \emptyset$ 
   $res \leftarrow \perp$ 
  While  $\pi \neq \emptyset$ 
     $r \leftarrow \pi.pop()$ 
     $res \leftarrow r.action.action()$ 
    If  $res = \perp$ 
       $r.action.repair()$ 
      break
    Else
       $E \leftarrow E \cup \{r\}$ 
    End if
  End while
  update_frequencies( $E$ ,  $res$ )
  ReturnReturn  $res$ 

```

Please note that if one would like to use a sliding window, she or he would also have to store data about whether  $R_i$  is element of  $\pi$  for all rules  $R_i$  and an executed plan  $\pi$  (i.e., the rows in  $A$ ) in a FIFO buffer. When we learn data about a new  $\pi_n$ , we might not only have to accommodate these new data, but some old  $\pi_o$  might fall out of the window so that we have to remove its influence. Via implementing a corresponding FIFO buffer, we could accomplish this easily.

With our definition of a plan’s reliability, we obviously search for a plan  $\pi$  with the lowest value  $r(\pi)$ .

Consequently, we do not reason about plans of minimal length, but desire plans of minimal costs in terms of  $r(\pi)$  which directly relates to the risk of a plan’s failing.

If there are multiple plans with the same “optimal”  $r(\pi)$ , we select the one created first. In principle, in such a case also heuristics that choose a plan of minimal length, or one such that the contained rules’ maximum reliability value is the lowest could be adopted.

In this context, it is also important to note that our concept is orthogonal to the incorporated planning algorithm/concept - as long as one can use our simple risk-related cost function drawing on results from light-weight diagnostics in that algorithm. So whether a derived plan  $\pi$  is globally or locally optimal (consider a complete vs. greedy search) depends on the incorporated planning algorithm. Consequently, the planning stage is not in the primary focus of our presentation, but we focus on (a) the exploitation of diagnostic data that describe the reliability of a system’s actions as well as on (b) how to exploit such data in the planning stage of a corresponding engine via a specific cost function in the form of a plan’s reliability (see Def. 7). While completeness and soundness in terms of finding a plan optimal in the context of the chosen heuristic (the suspiciousness coefficient like Ochiai) thus depends on the planning algorithm, we can easily deduce the complexity of our computations.

**Theorem 1** Computing a rule’s reliability coefficient is done in constant time, so that computing all of them is linear in the number of rules. Computing a plan  $\pi$ ’s reliability is linear in the length of its rule sequence.

**Proof:** Since we keep track of a component’s frequency values and only have to update them via simple additions and subtractions (the latter only in case of a sliding window), their computation and that of the chosen suspicious metric like Ochiai is in constant time. We do this for each rule so that we are linear in their number for the entire computation. Computing a plan  $\pi$ ’s reliability means summing up the contained rules’ values (see Def. 7) so that we’re linear in the length of  $\pi$ ’s sequence.

When implementing our concept in practice, there are some aspects that we have to consider in relation to the similarity coefficients, though. For instance, after a cold start, insufficient data would result in a

division by zero or a value of zero when computing the coefficients. Since we are using the values directly as reliability measures, and therefore for planning, in our proof-of-concept implementation, we assign small values (0.00001) as a rule’s reliability measure in such cases. This follows the idea that after a cold start, we assume that the rules are rather healthy than faulty.

As we stated in the previous section, PDDL is one of the most used planning description languages. To make our approach as universal as possible, we wanted to show that our approach can be used in combination with PDDL, and therefore can be integrated into every system that uses PDDL. However, using PDDL for our approach is not straight forward. We can see from Listing 1, that actions in PDDL are abstract actions still requiring concrete parameters when we want to execute that action. Because our approach uses the feedback from the actual execution of an action, it is a good idea to also reason about concrete actions and the current state of the world. Creating concrete actions from abstract actions is called grounding in different research areas. For our implementation of grounding, we have two important requirements. (1) The calculation of abstract actions to concrete actions has to be dynamic. (2) we have to be able to reuse already discovered concrete actions in the next plan.

To fulfill those two requirements, we implemented a reachability-based algorithm for computing concrete actions. Starting from the initial state, we (1) calculate which concrete actions could be executed from this state. For each such concrete action we (2) calculate how it transforms the state. If the state is a previously encountered state, we (3a) link the state to the already encountered state. If this is a new state, we (3b) continue at (1). The result is a graph where the nodes represent all different states of the world, and the edges represent their transformation through concrete actions. We can now use the reliability score of the concrete actions as edge weights and use traditional pathfinding algorithms to generate a plan. In our implementation, we use Dijkstra’s algorithm. Retaining this graph throughout different plans enables us to rediscover already used concrete actions.

It is important to note that with this approach, we reason with the maximum level of information about the concrete execution, i.e., the state of the world before the execution and the concrete action. This, of course, is not the only option. For example, reasoning only about abstract actions would also be possible, e.g., if we know the environment does not influence the result. Nevertheless, using our SFL approach to represent the confidence in an action’s success would still be applicable.

**Listing 1:** PDDL Domain for all examples

```
(define (domain robot-strips)
(:predicates (at ?r)(connected ?r1 ?r2)
(holding ?i)(itemat ?i ?r)(putlocation ?r))
(:action move
:parameters (?from ?to)
:precondition (and (at ?from)
(connected ?from ?to))
:effect (and (not (at ?from))
(at ?to)))
(:action pickup
:parameters (?room ?item)
:precondition (and (itemat ?item ?room)
(at ?room))
:effect (and (holding ?item)
(not (itemat ?item ?room))))
(:action put
:parameters (?room ?item)
:precondition (and (putlocation ?room)
(at ?room)
(holding ?item))
```

```
:effect (and (itemat ?item ?room)
             (not (holding ?item))))
```

**Listing 2:** PDDL Problem for “Shelves a priori”

```
(define (problem strips-robot)
  (:domain robot-strips)
  (:objects room_0_0 room_0_1 [...] item )
  (:init (at room_0_0)
         (putlocation room_0_0)
         (connected room_0_0 room_0_1)
         (connected room_0_1 room_0_0)
         [...])
  )
  (:goal (and )))
```

## Experiments

In order to evaluate our approach, we chose the scenario of a warehouse where an agent equipped with our intelligent reasoning system (further called intelligent agent) has to fetch randomly placed items while avoiding further agents. For representing such scenarios, we use a grid world consisting of discrete cells like illustrated in Figures 2, 3. Our intelligent agent can move around in its world via single actions for moving to the north, east, south, or west from its current cell. That is, if it is not blocked by a wall or another agent occupying the target cell (then it would remain in the current cell). Technically, we implemented our scenario as an extension to the OpenAI gym environment gym-maze (see <https://github.com/MattChanTK/gym-maze>).

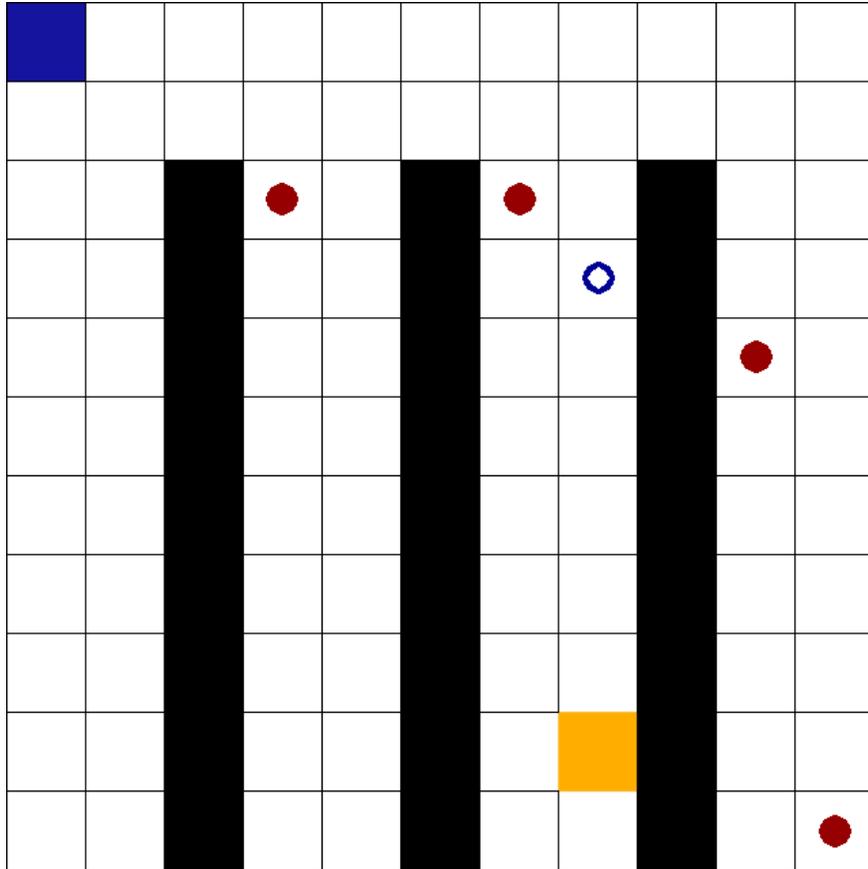


Figure 2: Warehouse example setup with shelves. The intelligent agent (blue circle) starts at 0/0 (blue square) and fetches items (orange square). Other agents are indicated by red dots.

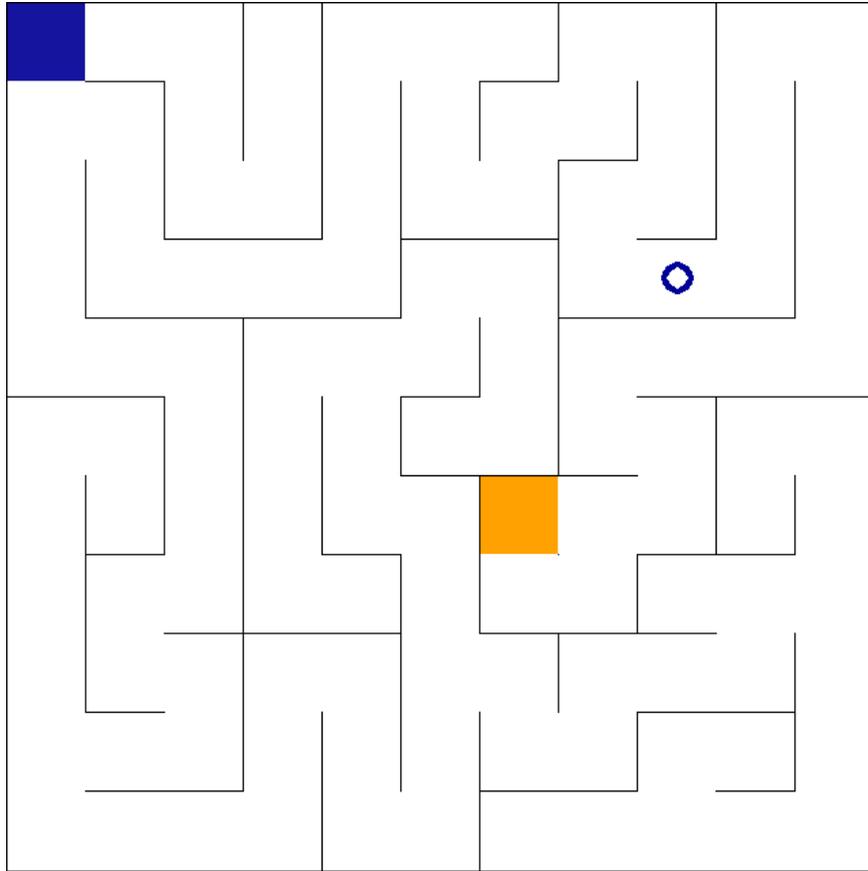


Figure 3: Warehouse example setup with a maze structure. The intelligent agent (blue circle) starts at 0/0 (blue square) and fetches items (orange square).

The intelligent agent knew about the warehouse domain and scenario from PDDL descriptions (Listing 1 and 2). In some experiments, also all walls were described in these files, in some the intelligent agent had to learn their location via the feedback of a failed move. Please note that it did not learn whether a move was blocked by an agent or a wall, so that these other agents add noise to the observed learning data.

All experiments use the same *PDDL domain file* shown in Listing 1 where we describe three actions. First, *move* enables the intelligent agent to move as described above. Second, *pickup* allows it to pick items up (if it is in the same cell), and third, via *put* it can put an item down in a put location.

The basic structure of all our *PDDL problem files* is shown in Listing 2}. Depending on the example configuration, specific atoms, i.e., (*connected room\_X1\_Y1 room\_X2\_Y2*), might be omitted from the initial state such as to indicate that there is a wall between the two cells. For each fetch task, the intelligent agent receives a new random item location, which is added to the PDDL initial state on-the-fly, and a PDDL goal to bring the item to the put location. The intelligent agent then starts at (0/0) in the grid, has to go to the item’s location, pick it up, go back to the put location (0/0), and, finally, deliver the item by putting it down and thus fulfilling the PDDL goal. If it is necessary to re-plan the intelligent agent will start from its current location.

We used different warehouse sizes (5x5, 8x8, or 11x11 cells) and numbers of other agents (0,1,4) in our experiments. The last parameter of a configuration is the setup of the experiment in terms of walls and an agent’s a priori knowledge of them as described below. When conducting the experiments, we investigated 100 different random fetch sequences for a specific configuration and reported the average values. Each such fetch sequence consisted of 100 fetch tasks. We will also show the average performance for each of them (over the 100 runs) such as to investigate the performance increase experienced. Please note that after finishing a fetch sequence, the learned knowledge was discarded. For reason of space, we report only on a few selected configurations in this section. The results for all configurations and the code for the experiments are available on GitHub (see <https://github.com/martinzimmermann/RBL-test-programs/releases/tag/CPS-RTSA>).

**Shelves a priori:** The grid world contains shelf cells that the agent cannot enter and around which two normal cells are placed (see Fig. 2). Items can only be located next to a shelf. For this setup, the PDDL Problem contains the shelves’ location (Lst. 2) so that an agent can move around efficiently. The challenge of this setup is that multiple agents are operating in the same warehouse. The intelligent agent knows nothing about their locations and can only learn about them by colliding with them. Still, the intelligent agent is expected to fetch items efficiently.

**Shelves a posteriori:** The setup is similar to the previous one, but the PDDL problem file does not contain data about the shelves. Thus, these data will be learned by the intelligent agent via the move actions’ reliability for the neighboring cells such as to be able to move around in the grid efficiently (while still having to avoid other agents).

**Maze:** The third setup is a randomly generated maze (see Fig. 3), where the intelligent agent does not know the layout of the maze but has to learn it through colliding with walls. Please note that since corridors are only one cell wide, it is not possible to bypass other agents. Thus there are no other agents in this setup.

## Experimental Results

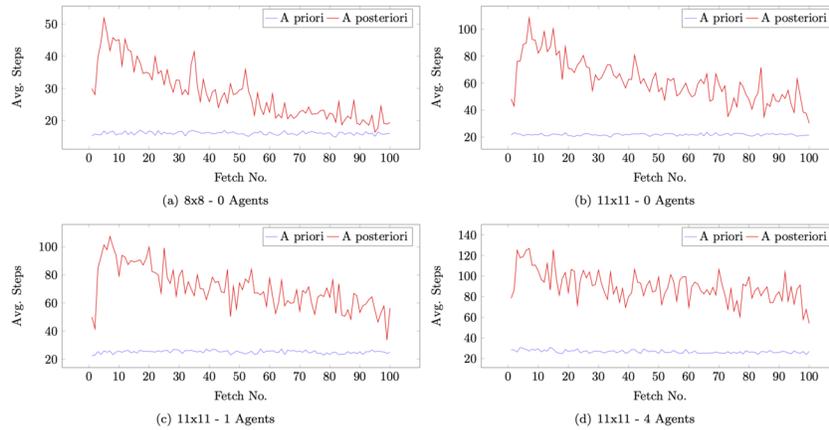


Figure 4: Results of selected configurations. The x-axis shows the specific fetch of the fetch sequence. The y-axis shows the average number of total steps needed for this fetch over all 100 fetch sequences. We can see that for configurations with 0 Agents the a posteriori avg. steps converge toward the a priori avg. steps fast (Subfigures a and b). However as the number of agents increases this gets slowed down (Subfigures c and d).

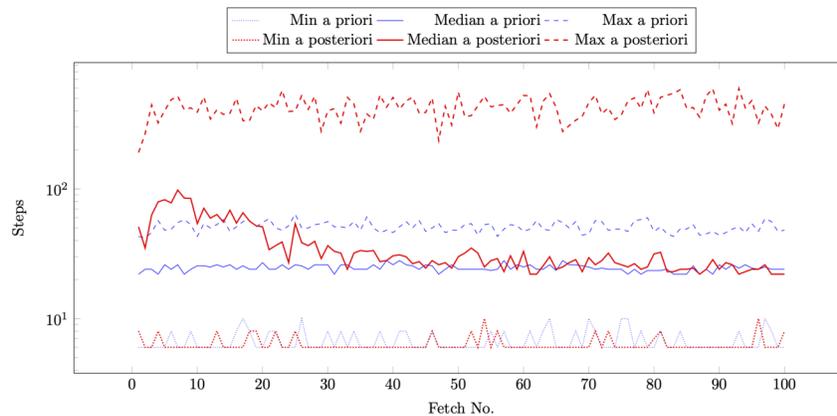


Figure 5: Min, median, and max steps over 100 runs for the 11x11 grid - 1 Agent configurations. The x-axis shows the specific fetch in the fetch sequence. The y-axis shows steps needed per fetch on a log scale. The graph shows that the minimum number of steps per fetch is similar for a priori and a posteriori, the median number of steps per fetch for a posteriori is in the beginning worse than a priori but converges toward the median of a priori over time. The max number of steps is high for a posteriori, but for a priori max remains low and stable over time.

In Fig. 4, we report the average number of steps needed per item fetch over all 100 runs. We can see that for all given configurations, the intelligent agent in the a posteriori setup needs fewer steps over time. This confirms our hypothesis that with our SFL approach, we can use the reliability measurement to enhance planning. However, we can see that this task gets more difficult when adding more agents, as this generates more random noise. The random noise makes it harder for the agent to distinguish between temporary

failures (i.e., other agents) and permanent faults (i.e., shelves).

For a priori, we can not confirm such an improvement. This is no surprise, as for this setup, the only unknown information about the world is the movement of the other agents. The movement is random and not learnable by our intelligent agent, as random behavior is, in general, not learnable. The small variance of steps needed can be explained by the random generation of the item locations. In Fig. 5, we can see that also the minimum, maximum, and median steps needed for a priori stay consistent over time.

One of our main focus points was also to compare the a priori and a posteriori setups. Fig. 4 shows that the a posteriori setup performs worse than the a priori setup. However, the problem of solving the a posteriori setup is much harder. First, it consists of more possible actions (moves through shelves are also considered during planning for a posteriori), and second, much information about the world, i.e., the location of the shelves, is unknown to the intelligent agent. It is remarkable that for the 8x8 a posteriori configuration with 0 agents a similar performance as a priori is reached after only around 100 fetches (Fig. 4a). This could be due to the reason that using 0 agents makes the scenario static, although still not known by the agent. For the other configurations, we also see a strong trend toward the performance of the a priori configurations. However, in our experiments, they never reach the same performance. It is not clear if just more fetches, meaning more training data, are needed to learn to distinguish between shelves and agents, or if they will never converge toward the a priori performance. To answer this, further experiments with longer fetch sequences are necessary.

During our investigation, we could not yet explain why, for most configurations, the performance of the first few fetches gets significantly worse before the performance gets better again. The only connection we could draw was that we sometimes saw similar behavior while performing reinforcement learning in a different domain. Further research is necessary to find the root cause of this behavior. However, this was not a major concern for us, as for all configurations, in the end we performed better than the first fetch.

In Fig. 5, we show the maximum, median, and minimum number of steps required per fetch over the 100 runs. Interestingly, the median, similar to the average, of a posteriori converges toward the median of a priori. The maximum, on the other hand, does not get smaller over time. The reason for this could be that the world is not sufficiently explored to calculate reliable plans for *all* locations after only 100 fetches. With roughly 100 cells where an item could be placed, this seems plausible. Similarly, there is a high chance that an item location is close to the start when considering 100 runs. Having an item close to the start would explain the very stable minimum for both a posteriori and a priori.

In Tab. ??, we compared average total steps, total plans, and total runtime summed over a whole run, i.e., 100 fetches, for the 11x11 grid a priori and a posteriori configuration. For the configuration a priori - 0 agents, the intelligent agent only needs 100 plans. That means no replanning was necessary. This was expected as everything about the world is known. As the number of agents increases, the number of replans increases. More replans are needed because there is a higher chance for the intelligent agent to collide with another agent. For the a posteriori setup, we see a similar increase in total plans and steps between 0 agents and 4 agents. But, the difference between the a priori and a posteriori is significant. One explanation for the difference could be that, for a priori, the information about the shelves is known, and the intelligent agent only collides with agents. In contrast, for a posteriori the intelligent agent, in the beginning, collides mostly with shelves. Because of lack of knowledge, the agent generates a plan that bypasses the shelf by just one block. Mostly, a shelf is right next to another shelf. This, in turn, leads to another collision with only a single step taken. The average steps per plan also support this explanation. For a priori, this is, on average, 15 steps per plan, and for a posteriori, this is, on average, only two steps per plan.

From Tab. ?? we can see that the runtime mostly correlates to the number of plans needed. Although we implemented some improvements for the RBL planning algorithm, it is still by far the largest runtime bottleneck. Only a fraction of the planning time is used for the SFL reliability calculation, which only increases with the number of possible actions for which the reliability should be calculated, as we proved in a previous section. Combining the SFL approach with a better planning algorithm would for sure result in

better runtime.

Configuration	Total Steps	Total Plans	Total Runtime
a prio - 0 A.	2184.72	100.0	38.74s
a prio - 1 A.	2507.48	155.56	50.63s
a prio - 4 A.	2669.05	349.95	95.69s
a post - 0 A.	6191.55	3111.57	3480.39s
a post - 1 A.	7072.04	3951.76	4374.32s
a post - 4 A.	9002.25	5798.05	5752.96s

Table 2: Results of the 11x11 configurations. The number is always the average over all 100 runs. Runtime was taken on a PC running Kubuntu 18.04.4 with an Intel(R) Core(TM) i7-6850K CPU@3.60GHz and 32GB RAM.

For the Maze, there were improvements similar to the shelves a posteriori configurations. This shows that we can also improve our plans in very complex environments.

## Comparison with other Approaches

In this section, we will compare RBL to other approaches in the literature. Because the different approaches excel in different areas, and no other approach focuses on the same topic as RBL, we performed a qualitative comparison between the approaches. First, we describe the other approaches and compare them with RBL. Afterwards, we compare the approaches all together considering seven different characteristics.

- **RL:** Reinforcement learning(RL) gained huge recognition as being able to produce agents that can dynamically act in diverse environments. Typically RL solves Markov decision problems where the agent tries to maximize the cumulative reward. Usually, this is done by first training an agent on a huge amount of sample data and then deploying it to a real situation. Although it is possible to train an RL agent also during operation, this is usually not done as the agent’s learning performance greatly depends on the exploration vs. exploitation tradeoff. The biggest difference between RL and RBL is that RL does not use any a priori knowledge and, therefore, always starts learning from scratch. With the exploitation of the knowledge given to RBL, it is possible to deploy RBL to a real situation directly and solve the problem right from the start. Even in the case that situations were unforeseen in the a priori knowledge, RBL can learn to circumvent these and stay operational. On the other hand, RL agents can achieve much better performance than RBL on a given task. During training, RL agents also take exploratory actions, which lead them to acquire new knowledge about the world. RBL only takes exploratory actions when a failure occurs and only in the scope of the knowledge provided.
- **POND:** POND provides an interesting approach to solve partial-observable and non-deterministic planning problems. It combines different search techniques and heuristics and switches between them dynamically. Furthermore, it employs a base representation for the problem from which other representations can be calculated, which are then used for the search or heuristic. Compared to RBL, it probably performs better when everything about the problem is known at the start time. However, POND lacks the capability to use feedback from the execution to reevaluate its plans, and it is also not able to learn new information. (Bryce, 2006)
- **FF-Replan:** FF-Replan was the winner of the 2004 International Probabilistic Planning Competition. It achieves this by first constructing a determinist planning problem out of the probabilistic planning problem and replanning when it encounters encounters encounters encounters encounters a state that differs from the expected one. The conversion from a probabilistic plan to a determinist plan is either done by single-outcome or all-outcome. Single-outcome chooses a single action among many probable ones depending on a heuristic, and all-outcome creates a separate action for each. FF-replan in all-outcome mode is quite similar to our approach. In RBL, we are not concerned with probabilistic

actions per se. However, we permit each action to fail without knowing before which action and when the action will fail. We could emulate such behavior in FF-replan by adding a fail outcome to every action. However, FF-replan would still not learn from failures like RBL and would probably get stuck as soon as reality would not conform to its knowledge (e.g., a wall is in reality where there is non in the model). (Yoon et al., 2007)

- **PRM-RL:** PRM-RL combines Probabilistic Roadmaps (PRM) and RL. It first trains an RL agent via Monte Carlo selection on a similar environment as the target environment to get an agent that can successfully move in the environment. After this step, the agent is deployed to the target environment, and the PRM builder creates a PRM based on a uniform sampling of the agent’s movement. Only collision-free point-to-point navigation is retained in the PRM. After these two training steps, PRM-RL can successfully navigate the target environment. Although PRM-RL currently lacks the ability to learn in operation like RBL, it is not hard to imagine that the PRM builder could also be run as soon as the PRM model differs from the environment and therefore signals that a fault occurred. A benefit of PRM-RL is that it does not need a priori knowledge, and it infers everything from training. However, the RL agent and the PRM builder need this training to function properly, compared to RBL, which can be used without training at all. (Faust et al., 2017)

For the characteristics, we selected “Needs model,” “Needs Probabilities,” “Needs Training,” “Learns during operation,” “Performs Exploration,” “Failure resilient,” and “Guarantees.” We think these are the most important characteristics when someone wants to decide which approach he should use. Our results are presented in Table 3, and a detailed description of the characteristics can be found below.

Ap- proach	Needs model	Needs Probabilities	Needs Training	Learns during operation	Performs Exploration	Failure resilient	Guar- antees
RBL	x	-	-	x	-	x	x
RL	-	-	x	-	x	x	-
POND	x	x	-	-	x	-	-
FF- Replan	x	x / -	-	-	-	x	x
PRM- RL	-	-	x	-	x	-	-

Table 3: Qualitative comparison of different approaches. Note that although FF-Replan accepts probabilistic PDDL, the probabilities are only used in one variant of the approach.

**Needs model:** Describes if the approach needs a model to be usable. This can be in the form of PDDL, PDDL-like, or other non-formal information. Usually, approaches that have this model available perform better than others as they do not have to learn a model first. However, it is not always easy to get a model. Here “x” means the approach needs this information, and “-” means the approach does not need the information.

**Needs Probabilities:** Describes if the approach needs to know probabilities of the non-deterministic actions to function. Like the approach before, this information is often not easy to obtain or even impossible in a dynamic scenario. Here “x” means the approach needs this information, and “-” means the approach does not need the information.

**Needs Training:** Describes if and how much training the approach needs before it can be used. “x”, means training is required before the approach can be used, e.g., in a simulation, “-” the approach either does not need training or will learn during operation.

**Learns during operation:** Describes if the approach can learn during operation. “x” means that the approach will learn during operation, “-” means the approach is fixed during operation.

**Performs Exploration:** Describes if the approach can perform exploratory actions. This is useful if only partial information is available, and for example, new actions are tried or actions for knowledge gain are performed.

**Failure resilient:** Describes if the approach can handle unforeseen circumstances. For example, if the approach is able to adapt if a fault occurs during operation. “x” means the approach can deal with unforeseen circumstances, “-” means the approach will fail if an unforeseen circumstance occurs.

**Guarantees:** Describes if the approach gives some guarantees.

*FF-Replan:* Because FF-Replan only selects actions with a non-zero percent chance of leading to a goal, FF-Replan is guaranteed to reach the goal eventually if there are no dead ends and the environment is equivalent to the provided model.

*RBL:* Similar to FF-Replan also RBL takes only actions that have a chance to lead to a goal. However, because RBL also updates the reliability of the actions, the environment does not have to be equivalent to the provided model. Therefore, we can guarantee that an agent with RBL will eventually reach the goal if there are no dead ends and there is a possible action sequence possible in the model that would lead to the goal. Meaning, as long as there are redundancies in the model of which not all are blocked.

## Related Work

Nilsson presented in his work (Nilsson, 1994) with Teleo-Reactive programs a formalism for action sequences an agent can take to reach goals in uncertain environments. Teleo-Reactive programs are an ordered list of actions with preconditions. The first action of this list which preconditions are met, is executed indefinitely. Through clever construction of Teleo-Reactive programs, an agent can then deal with uncertainty.

Krenn and Wotawa (Krenn and Wotawa, 2009) proposed a similar approach. Instead of just using the first rule of the list, the rules’ selection frequency could be dynamically updated during operation. The rules’ selection frequency was in this approach based on biological processes of DNA transcription. Furthermore, rules do not only have precondition, but preconditions and postconditions.

Our approach builds upon this line of research. Our approach extends this further with two main contributions. First, we created an interface to PDDL, enabling the approach to apply to a wide range of already existing models. Second, instead of the biological inspiration, we use SFL, which was already proven to be successful at detecting faulty components in software testing, which is more akin to the problem of distinguishing faulty from non-faulty actions.

Another related line of research is replanning and plan repair. (Chien et al., 2000) and (Aschwanden et al., 2006) use feedback from the execution to change the current beliefs of the system. (Wilkins, 1988) plans multiple plans and chooses a new one when one fails. (Draper et al., 2013) handles uncertainty by incorporating conditions into the plan. Georgeff follows in (Georgeff and Lansky, 1987) a similar approach and uses partial planning and delayed decisions.

All those approaches have in common that, in case of a failure, they have to either change the beliefs, locate the fault, or gather additional information about the environment to choose a (new) plan. In contrast, our approach does not have to locate the fault but rather uses statistics to locate the plan’s reliable actions. In particular, the calculation with SFL is very inexpensive and can be adapted for different approaches. Even extending former mentioned approaches should be possible as long as we can integrate a cost function into the planning algorithm. To the best of our knowledge, we do not know of a similar approach for replanning.

Finally, also Reinforcement Learning (RL) is related to our approach. RL, as previously mentioned, is capable of solving a wide variety of tasks. Modern RL consists mostly of two intertwined research directions. First, try-and-error learning inspired by (Minsky, 1961). Second, Optimal Control, which has its beginnings in the 1950s, mostly by Bellman. In (Bellman, 1954) he proposed an approach, Dynamic Programming, which could solve optimal control programs. However, this approach did not scale well for higher dimensional spaces. Sutton and Barto took inspiration from these two approaches and mixed it with temporal-difference to create the modern RL (Sutton and Barto, 1981). A lot of improvements and demonstrations were made over the last few years (Silver et al., 2017; Berner et al., 2019).

However, a core problem still remains that RL agents usually need many training samples till they are operational. Also, little research about how to bootstrap an agent with models like PDDL was done. Our approach uses a model as a bootstrap process and further refines the reliability of the available actions by try-and-error. This allows our approach to perform reasonably well right from the start without first learning how to operate in the environment. Another problem for RL is constructing a reward function, which is quite tricky to get right (Clark and Amodei, 2016). In our approach, no reward function has to be defined as we only learn from action failures.

## Conclusion

We showed how to adopt SFL for a live setting in order to generate a metric for catching our actions/rules' reliability or healthiness. We used the computed similarity coefficient values directly for selecting future rule sequences that are most likely to succeed in achieving our goals - following the idea that a sequence's risk of failing directly corresponds to the sum of the individual action's risk's of failing as expressed by our SFL metric. We showed how to easily compute these values dynamically (in constant time for a single rule) and how to adopt a sliding window if desired for a highly dynamic environment. Combining SFL diagnostics with a planning and execution environment like RBL enabled us to foster intelligent behavior taking the constantly observed reliability of a system's actions into account. Our experiments showed that we indeed can profit from our learning about the action's reliability. Although neither using feedback from a plan's execution to improve planning, nor using SFL for rule-bases are novel in general, combining both and adopting our concept in our context are indeed novel contributions that lead to attractive results and are (a) easy to adopt and (b) easy to compute such that it fits also applications in embedded cyber-physical systems where resources might come at a premium. Furthermore, we gave a qualitative comparison between RBL and other related approaches. With the help of this comparison the reader can deduce different trade-offs of the approaches and select the appropriate approach for his scenario.

While not reported in detail for our experiments, please note that we experimented with several similarity coefficients and found Jaccard to work best for our configurations. Future experiments investigating also sliding windows and longer learning phases will have to confirm such first trends though - also in the context of multiple scenario domains.

Further room for improvement is the tuning of the concrete values that are added to the spectrum (currently only  $1$  and  $0$ ). From reinforcement learning, we see that discounting rewards based on their temporal ordering positively influences the learning rate of the system. A similar approach could also be taken for the values in SFL, including exploration stages with specific strategies. That is, entirely unlimited exploration could be exploited to gather broader knowledge at the cost of performance in the tasks themselves, but also limiting the exploration to plans that deviate in performance only within some boundary to the optimal one could provide a more limited but still more educated picture. In such future research, it will also be interesting to consider effects from temporal considerations when associating the blame of a plan's failure to individual actions, or considering previous executions from the less important in the SFL spectrum as recent ones.

## Acknowledgements

The financial support by the Austrian Federal Ministry for Digital and Economic Affairs and the National Foundation for Research, Technology, and Development is gratefully acknowledged.

## Conflict of interest

The authors declare not to have any financial or commercial conflicts of interest.

## References

- Rui Abreu, Peter Zoetewij, Rob Golsteijn, and Arjan J.C. van Gemund. A practical evaluation of spectrum-based fault localization. *Journal of Systems and Software*, 82, 2009a.
- Rui Abreu, Peter Zoetewij, and Arjan J. C. van Gemund. Spectrum-Based Multiple Fault Localization. In *ASE*, pages 88–99, 2009b.
- Pascal Aschwanden, Vijay Baskaran, Sara Bernardini, Chuck Fry, Maria Moreno, Nicola Muscettola, Chris Plaunt, David Rijsman, and Paul Tompkins. Model-Unified Planning and Execution for Distributed Autonomous System Control. In *Spacecraft Autonomy: Using AI to Expand Human Space Exploration, Papers from the 2006 AAAI Fall Symposium*, 2006.
- Richard Bellman. The theory of dynamic programming. *Bulletin of the American Mathematical Society*, 60, 1954.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Christopher Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Pondé de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with Large Scale Deep Reinforcement Learning. *CoRR*, abs/1912.06680, 2019. URL <http://arxiv.org/abs/1912.06680>.
- Avrim L. Blum and Merrick L. Furst. Fast planning through planning graph analysis. *Artificial Intelligence*, 90(1):281 – 300, 1997. ISSN 0004-3702. doi: [https://doi.org/10.1016/S0004-3702\(96\)00047-1](https://doi.org/10.1016/S0004-3702(96)00047-1). URL <http://www.sciencedirect.com/science/article/pii/S0004370296000471>.
- Craig Boutilier, Raymond Reiter, Mikhail Soutchanski, and Sebastian Thrun. Decision-Theoretic, High-Level Agent Programming in the Situation Calculus. In *17th National Conference on Artificial Intelligence and 12th Conference on Innovative Applications of Artificial Intelligence*, pages 355–362, 2000. URL <http://www.aaai.org/Library/AAAI/2000/aaai00-055.php>.
- Daniel Bryce. POND: The partially-observable and non-deterministic planner. *ICAPS 2006*, page 58, 2006.
- Steve Chien, Russell Knight, Andre Stechert, Rob Sherwood, and Gregg Rabideau. Using Iterative Repair to Improve the Responsiveness of Planning and Scheduling. In *Proceedings of the Fifth International Conference on Artificial Intelligence Planning Systems*, 2000.
- Jack Clark and Dario Amodei. Faulty Reward Functions in the Wild, 2016. URL <https://openai.com/blog/faulty-reward-functions/>.
- D. Connell and H. M. La. Dynamic path planning and replanning for mobile robots using RRT. In *2017 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 1429–1434, 2017. doi: 10.1109/SMC.2017.8122814.

- Johan de Kleer and Brian C. Williams. Diagnosing Multiple Faults. *Artificial Intelligence*, 32(1):97–130, 1987.
- Denise L. Draper, Steve Hanks, and Daniel S. Weld. A Probabilistic Model of Action for Least-Commitment Planning with Information Gather. *CoRR*, abs/1302.6801, 2013. URL <http://arxiv.org/abs/1302.6801>.
- Aleksandra Faust, Oscar Ramirez, Marek Fiser, Kenneth Oslund, Anthony G. Francis, James Davidson, and Lydia Tapia. PRM-RL: Long-range Robotic Navigation Tasks by Combining Reinforcement Learning and Sampling-based Planning. *CoRR*, 2017. URL <http://arxiv.org/abs/1710.03937>.
- Richard E. Fikes and Nils J. Nilsson. Strips: A new approach to the application of theorem proving to problem solving. *Artificial Intelligence*, 2(3):189 – 208, 1971. ISSN 0004-3702. doi: [https://doi.org/10.1016/0004-3702\(71\)90010-5](https://doi.org/10.1016/0004-3702(71)90010-5). URL <http://www.sciencedirect.com/science/article/pii/0004370271900105>.
- Michael P. Georgeff and Amy L. Lansky. Reactive Reasoning and Planning. In *Proceedings of the Sixth National Conference on Artificial Intelligence - Volume 2*, AAAI’87, pages 677–682. AAAI Press, 1987. ISBN 0-934613-42-7. URL <http://dl.acm.org/citation.cfm?id=1863766.1863818>.
- M. Ghallab, A. Howe, C. Knoblock, D. Mcdermott, A. Ram, M. Veloso, D. Weld, and D. Wilkins. PDDL—The Planning Domain Definition Language, 1998. URL <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.37.212>.
- Stephan Gspandl, Ingo Pill, Michael Reip, Gerald Steinbauer, and Alexander Ferrein. Belief Management for High-Level Robot Programs. In *IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22 2011*, pages 900–905, 2011. URL <https://doi.org/10.5591/978-1-57735-516-8/IJCAI11-156>.
- James A. Jones and Mary Jean Harrold. Empirical evaluation of the tarantula automatic fault-localization technique. In *Proceedings ASE*, 2005.
- Henry Kautz and Bart Selman. Pushing the Envelope: Planning, Propositional Logic, and Stochastic Search. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence - Volume 2*, AAAI’96, pages 1194–1201. AAAI Press, 1996. ISBN 0-262-51091-X. URL <http://dl.acm.org/citation.cfm?id=1864519.1864564>.
- Willibald Krenn and Franz Wotawa. Intelligent, Fault Adaptive Control of Autonomous Systems. In Natividad Martinez Madrid and Ralf E.D. Seepold, editors, *Intelligent Technical Systems*, pages 175–188. Springer Netherlands, Dordrecht, 2009. ISBN 978-1-4020-9823-9. doi: 10.1007/978-1-4020-9823-9\_13. URL [https://doi.org/10.1007/978-1-4020-9823-9\\_13](https://doi.org/10.1007/978-1-4020-9823-9_13).
- M. Minsky. Steps toward Artificial Intelligence. *Proceedings of the IRE*, 49(1):8–30, 1961. doi: 10.1109/JRPROC.1961.287775.
- Nils J. Nilsson. Teleo-Reactive Programs for Agent Control. *J. Artif. Int. Res.*, 1, 1994.
- Ingo Pill and Franz Wotawa. Spectrum-Based Fault Localization for Logic-Based Reasoning. In *2018 IEEE International Symposium on Software Reliability Engineering Workshops, ISSRE Workshops, Memphis, TN, USA, October 15-18, 2018*, pages 192–199, 2018. URL <https://doi.org/10.1109/ISSREW.2018.00006>.
- Raymond Reiter. A Theory of Diagnosis from First Principles. *Artificial Intelligence*, 32(1):57–95, 1987.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy P. Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm. *CoRR*, abs/1712.01815, 2017. URL <http://arxiv.org/abs/1712.01815>.

- R. Sutton and A. Barto. Toward a modern theory of adaptive networks: expectation and prediction. *Psychological review*, 88 2:135–70, 1981.
- David E. Wilkins. *Practical Planning: Extending the Classical AI Planning Paradigm*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1988. ISBN 0-934613-94-X.
- Sungwook Yoon, Alan Fern, and Robert Givan. FF-Replan: A Baseline for Probabilistic Planning. In *Proceedings of the Seventeenth International Conference on International Conference on Automated Planning and Scheduling*, 2007.
- Martin Zimmermann and Franz Wotawa. An adaptive system for autonomous driving. *Software Quality Journal*, 1 2020. ISSN 0963-9314. doi: 10.1007/s11219-020-09519-w.