



**HAL**  
open science

# Mutation-guided Metamorphic Testing of Optimality in AI Planning

Quentin Mazouni, Arnaud Gotlieb, Helge Spieker, Mathieu Acher, Benoit Combemale

## ► To cite this version:

Quentin Mazouni, Arnaud Gotlieb, Helge Spieker, Mathieu Acher, Benoit Combemale. Mutation-guided Metamorphic Testing of Optimality in AI Planning. *Journal of Software Testing, Verification and Reliability*, 2024, 5 (1), pp.e1898. <10.1002/stvr.1898>. <hal-05494710>

**HAL Id: hal-05494710**

**<https://hal.science/hal-05494710v1>**

Submitted on 5 Feb 2026

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

# Mutation-guided Metamorphic Testing of Optimality in AI Planning

Quentin Mazouni<sup>a</sup>, Arnaud Gotlieb<sup>a,\*</sup>, Helge Spieker<sup>a</sup>, Mathieu Acher<sup>b</sup>, Benoit Combemale<sup>b</sup>

<sup>a</sup>*Simula Research Laboratory, Oslo, Norway*

<sup>b</sup>*INRIA-IRISA-INSa, Rennes, France*

---

## Abstract

Autonomous systems such as space- or underwater-exploration robots or elderly people assistance robots often include an Artificial Intelligence (AI) planner as a component. Starting from the initial state of a system, an AI planner automatically generates sequential plans to reach final states that satisfy user-specified goals. Generating plans having a minimum number of intermediate steps or taking the least time to execute is usually strongly desired, as these plans exhibit minimal costs. Unfortunately, testing if an AI planner generates optimal plans is almost impossible because the expected cost of these plans is usually unknown. Based on mutation adequacy test suite selection, this article proposes a novel metamorphic testing framework for detecting the lack of optimality in AI planners. The general idea is to perform a systematic but non-exhaustive state space exploration from the initial state and to select mutant-adequate states to instantiate new planning tasks as follow-up test cases. We then check a metamorphic relation between the automatically generated solutions of the AI planner for these new test cases and the cost of the initial plan. We implemented this metamorphic testing framework in a tool called MORPHINPLAN. Our experimental evaluation shows that MORPHINPLAN can detect non-optimal behaviour in both mutated AI planners and off-the-shelf, configurable planners.

---

\*Corresponding author

*Email addresses:* [quentin@simula.no](mailto:quentin@simula.no) (Quentin Mazouni), [arnaud@simula.no](mailto:arnaud@simula.no) (Arnaud Gotlieb), [helge@simula.no](mailto:helge@simula.no) (Helge Spieker), [mathieu.acher@irisa.fr](mailto:mathieu.acher@irisa.fr) (Mathieu Acher), [benoit.combemale@irisa.fr](mailto:benoit.combemale@irisa.fr) (Benoit Combemale)

It also shows that our proposed mutation adequacy test selection strategy outperforms three alternative test generation and selection strategies, including both random state selection and random walks through the state space in terms of mutation scores.

*Keywords:* Metamorphic Testing, AI Planning, Mutation Testing

---

## 1. Introduction

### 1.1. Context

In many real-life autonomous systems that involve decision-making, such as space-exploration or underwater robots, elderly people assistance robots or automated production planning and industrial control systems, AI planning can be used to find a sequence of actions which moves a system from an initial state into a final state that satisfies a specified goal [20]. The recent development of competitive AI planners has opened an arena in Europe [39] and the world<sup>1</sup> for embedding more and more of these tools into complex autonomous systems.

As illustrated in Figure 1, AI planners take as input the definition of a *domain* which describes the predicates and actions of the considered environment and a *problem* which describes the initial state of the system and a goal to reach. As output, these planners return a plan, i.e., a sequence of instantiated actions with the objects of the problem, which moves the system from its initial state to a final state compliant with a user-specified objective. AI planners are expected to be sound and complete, which means that they shall always return a valid plan, i.e., a plan which correctly satisfies the assigned goal of the problem with domain definition, when there exists one, and they shall return “*no solution*” when there is none. Some AI planners can in addition return plans which optimize a given cost function. Common cost functions to minimize are, for example, the number of actions in the plan, the execution time of the plan, or the number of resources required to execute the plan. The language in which AI planning

---

<sup>1</sup>International Planning Competition 2023 - <https://ipc2023.github.io/>

23 problems are expressed, called *Planning Domain Definition Language (PDDL)*  
 24 [21], has been standardized to facilitate the comparison between AI planners  
 25 and foster the reuse of existing complex domain and problem definitions.

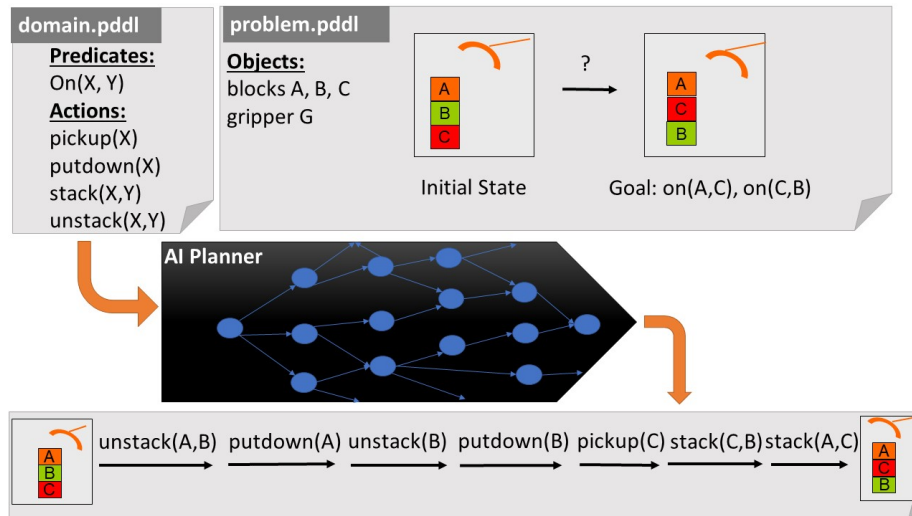


Figure 1: Solving an AI Planning problem in a nutshell.

## 26 1.2. Motivations

27 The development of AI planners in safety-critical missions involves strong  
 28 software engineering practices such as systematic validation and verification.  
 29 Many research works have addressed the problem of checking that AI planners  
 30 are sound and complete [37, 6, 4, 54, 18]. However, to the best of our knowledge,  
 31 the exact methods to determine if a returned plan is optimal or not do not  
 32 scale up. Even though this question can be less crucial than soundness and  
 33 completeness, some autonomous systems have to perform activities which span  
 34 over a long-term time horizon and they cannot afford to lose time or resources  
 35 by running unnecessary actions. For instance, communicating with a Mars rover  
 36 takes in between 8 to 42 minutes, and it can only be done a few times during the  
 37 Martian day, so the rover must compute optimal plans to save its energy and

38 communication resources. In production planning, computing optimal plans is  
39 often a crucial requirement to minimize the item’s delivery time.

40 The predominant approach for verifying the results of AI planners is based  
41 on formal verification [46, 41, 4, 45, 1]. Here, the main idea is to check formal  
42 temporal logic properties of the underlying transition system by using (bounded)  
43 model checking. Although this approach is appealing, it is demanding in terms  
44 of specifications and is currently focused on verifying only safety-critical proper-  
45 ties. Other non-functional property such as running time or memory consump-  
46 tion of the generated plan remains insufficiently explored in model checking of  
47 AI planners. Current methods for testing AI planners, which are less demand-  
48 ing than formal verification and also broadly used [27], are focused either on  
49 testing to which extent these planners deal with PDDL compliance [26, 38] or  
50 on testing their inputs domain and problem definition [40, 22, 32], i.e., ensuring  
51 that the definitions of actions do not contain contradiction which would pre-  
52 vent their application. Some works have also considered using AI planning to  
53 generate test cases for traditional software systems [25, 5]. But, despite the  
54 interest in verifying formally the correctness of generated plans [19] or testing  
55 domain/problem definition [22], there is a great challenge in checking whether  
56 generated plans are optimal or not. Besides detecting non-optimal plans which  
57 are unwanted, it is critical to determine if longer time or additional resources  
58 granted to a given AI planner could lead to less costly plans. This is crucial  
59 for safety- or business-critical missions with only limited time and resources to  
60 compute plans.

61 Testing AI planners to detect when non-optimal plans are generated is chal-  
62 lenging as, in general, no oracle is available for a given planning problem. Apart  
63 from simple domains and problems where analytical optimal solutions can be  
64 computed (e.g., Hanoi tower problems [17]), optimal solutions are unknown and  
65 computing those would require exploring the overall state space, which is infeas-  
66 ible in general. An alternative can be to use differential testing [11] by launching  
67 several distinct planners and comparing the costs of the generated plans. A  
68 difference between the costs would reveal a fault in at least one of the planners.

69 However, this approach is not suitable for autonomous systems when tests have  
70 to be performed online, as only limited computational resources are available.  
71 AI planners are very demanding tools in terms of memory as they have to ex-  
72 plore large state spaces, and it is often too costly to concurrently launch several  
73 planners on a single machine. Besides, AI planners for autonomous systems  
74 are usually meant to run on specific hardware and architecture, which does not  
75 facilitate the usage of multiple planners.

### 76 1.3. Contributions of the Article

77 In this article, we alleviate the problem by proposing mutation-guided meta-  
78 morphic testing for detecting non-optimal AI planners. *Metamorphic Testing* is  
79 a testing technique introduced in [8] to replace test oracle checking with *meta-*  
80 *morphic relations (MRs)* to assess the results of multiple program executions  
81 [43, 9]. Metamorphic Testing is used to generate so-called *follow-up test cases*  
82 [8, 10] using MRs, and it was successfully deployed to test a variety of complex  
83 software systems including ML systems [50, 51, 52, 47, 49], scientific software  
84 [53, 28], virtual reality applications [13] or autonomous vehicles [56, 14, 2]. Us-  
85 ing mutation-guided test case selection, we generate follow-up test cases (from  
86 reachable states of the state space) from an initial problem. Then, we check  
87 the cost of the initially generated plan and the costs returned by the solutions  
88 generated by the planner under test, from the follow-up test cases. By assuming  
89 the optimality of the AI planner, our metamorphic relation states that the over-  
90 all cost of the initial plan must be lower than or equal to the cost of the plan for  
91 any follow-up test case added to the cost for reaching the corresponding state.  
92 We evaluated our framework by comparing its ability to detect non-optimal be-  
93 haviour in two off-the-shelf, configurable planners, considered under 20 different  
94 configurations. Our mutation-adequacy strategy led to better detection of non-  
95 optimal plans and significantly lowered the number of required follow-up test  
96 cases, as compared with simpler baseline strategies including random genera-  
97 tion. However, its running time is also significantly higher than these simpler  
98 strategies and thus it is crucial to control its usage as discussed in the article.

99 The key contributions are:

- 100 • We present a mutation-guided metamorphic testing framework to detect  
101 non-optimal planning. It generates follow-up test cases from a set of states  
102 selected after an exploration of the planning problem’s state space. Viola-  
103 tions of the metamorphic relation let us detect non-optimal planning and  
104 thus reveal faults in the program implementation.
- 105 • We handle the follow-up test case generation in two steps: state space ex-  
106 ploration and state selection. The selection step scores the explored states  
107 with their ability to kill mutated AI planners. We choose these mutants  
108 by combining mutated functions leveraged in a classical  $A^*$  search. This  
109 approach (distinct from traditional mutation operators) lets us preserve  
110 the validity of the mutated programs’ outputs.
- 111 • We implemented this framework in a tool called MORPHINPLAN. To the  
112 best of our knowledge, it is the first tool of that kind, able to detect  
113 non-optimal plans returned by existing AI planners.
- 114 • Using mutation testing, we conduct a two-stage empirical evaluation of  
115 our test case generation method, for testing two off-the-shelf planning  
116 systems, including Fast Downward. We show that MORPHINPLAN can  
117 successfully generate test cases that detect non-optimal plans in mutated  
118 AI planners and existing configurations. The evaluation compares the  
119 mutation adequacy state selection strategy with three alternative state  
120 selection strategies used as comparison baselines and reveals an improved  
121 rate of fault-revealing follow-up test cases, thus improving the efficiency  
122 of testing.

#### 123 *1.4. Plan of the Article*

124 The article is organized as follows: Section 2 introduces AI planning and  
125 Metamorphic Testing with an illustrative example. Section 3 introduces the  
126 Metamorphic Testing framework. Section 4 describes the MORPHINPLAN tool

127 while Section 5 presents our empirical evaluation. Threats to the latter are  
128 treated in Section 6. Section 7 compares our work with other approaches. Even-  
129 tually, Section 8 concludes the article and draws some perspectives.

## 130 2. Background

### 131 2.1. AI Planning

132 Given an initial state and a goal to reach, AI planning aims at generating  
133 a plan from the initial state to a final state which satisfies the goal. A *plan*  
134 (Def.2) can be seen as a sequence of successive actions in a *state transition*  
135 *system* (Def.1).

136 **Definition 1** (State Transition System). *A State Transition System is a tuple*  
137  $\Sigma = (S, A, c, \gamma)$  *where*  $S$  *is a set of states,*  $A$  *is a set of actions,*  $c : A \rightarrow \mathbb{R}^+$   
138 *is a cost function,*  $\gamma : S \times A \rightarrow S$  *is a state transition function.*

139 In a given state  $s$ , only a limited number of actions can usually be applied  
140 and thus  $\gamma$  is only a partial function. When  $\gamma(s, a)$  is defined, we say that the  
141 action  $a$  is *applicable* in the state  $s$  and that it results in a new state  $s'$  with  
142 cost  $c(a)$ .

143 **Definition 2** (Plan). *A plan is a sequence of actions*

144  $\pi = [a_1, \dots, a_k]$  *s.t. each action of*  $\pi$  *is applicable from an initial state*  $s_0$  *of*  
145 *the system. Applying*  $\pi$  *to*  $s_0$  *leads to the following sequence of states:*  $s_1 =$   
146  $\gamma(s_0, a_1), s_2 = \gamma(s_1, a_2), \dots, s_k = \gamma(s_{k-1}, a_k)$ . *In that context,*  $s_k$  *is called the*  
147 *final state of the system after the application of*  $\pi$ . *Note that the size of*  $\pi$  *is*  $k$ ,  
148 *noted*  $|\pi|$  *and its cost is*  $\sum_{i=1}^k c(a_i)$ , *noted*  $c(\pi)$ .

149 A state  $s$  is reachable from an initial state  $s_0$  if there exists a plan  $\pi =$   
150  $[a_1, \dots, a_k]$  s.t.  $s = \gamma(\dots(\gamma(\gamma(s_0, a_1), a_2)\dots), a_k)$ . By extension,  $\gamma(s_0, \pi)$  denotes the  
151 state that results from applying  $\pi$  to  $s_0$ . In AI planning, there are several repre-  
152 sentations of states and actions. We adopt the so-called “logical representation”  
153 where each state and action modification is captured by logical propositions. All

154 examples of the article use PDDL syntax [21]. We assume some basic knowledge  
155 of PDDL, even though it is not mandatory to understand the article.

156 A state  $s$  is represented with a conjunction of literals over constant symbols  
157 (e.g.  $r$ ). Variables can be used to represent multiple symbols. Variables are  
158 noted with a question mark (e.g.,  $?r$ ).

159 **Example 1.** (*State*)

160 `(at r1 w)` `(loaded r1 i)` describes a state where a known rover  $r1$  stands at a  
161 waypoint  $w$  and holds an item  $i$ , while `(at ?r w)` is true iff there is a rover  
162 (unknown) standing at waypoint  $w$ .

163 An action  $a$  is a triple  $(head(a), pre(a), eff(a))$  where (i)  $head(a)$  is com-  
164 posed of a name and a list of variable parameters; (ii)  $pre(a)$  is a set of literals  
165 which must stand for the application of the action; (iii)  $eff(a)$  is a set of effects  
166 expressed with literals. These literals become true after the application of the  
167 action.

168 **Example 2.** (*Action*)

169 `(:action move :parameters (?r ?w1 ?w2)`  
170 `:precondition (at ?r ?w1)`  
171 `:effect (and (at ?r ?w2) (not (at ?r ?w1))))`

172 The `move` action takes as inputs a rover  $r$  and two waypoints  $w1, w2$ . Given that  
173  $r$  is at  $w1$ , the action moves the rover to  $w2$ , which is not anymore at  $w1$ . The  
174 use of variables makes the action generic for all rovers and waypoints.

175 Using the logical representation of states and actions, an AI planner aims at  
176 solving *planning problems*, defined as follows:

177 **Definition 3** (AI Planning Problem). Given  $P$  a finite set of propositions,  $A$  a  
178 finite set of actions, each of them being represented by the triple  $(head(a), pre(a), eff(a))$   
179 where  $pre(a)$  and  $eff(a)$  are subsets of  $P$ , a planning problem is defined as a  
180 tuple  $(P, A, s_0, obj)$  where, in addition to  $P$  and  $A$ , an initial state  $s_0 \subseteq P$  is  
181 defined and an objective  $obj \subseteq P$ , composed of conjunction of literals is given.

182 In PDDL, AI planning problems are materialized by two files, namely a  
 183 *domain* file which contains the definition of predicates and actions and a *problem*  
 184 file which contains the definition of the initial state and the objective. The  
 185 following domain example shows the actions that can be performed by a very  
 186 simple rover.

187 **Example 3.** (*Rover Domain Definition*)

```
188 (define (domain Rover)
189 (:predicates (rover ?r) (free ?r) (wp ?w) (item ?i) (at ?r ?w)
190 (loaded ?r ?i))
191
192 (:action load :parameters (?r ?w ?i)
193 :precondition (and (rover ?r) (wp ?w) (item ?i) (at ?r ?w) (at ?i ?w)
194 (free ?r))
195 :effect (and (loaded ?r ?i) (not (free ?r)) (not (at ?i ?w))))
196
197 (:action unload :parameters (?r ?w ?i)
198 :precondition (and (rover ?r) (wp ?w) (item ?i) (at ?r ?w) (loaded ?r ?i))
199 :effect (and (free ?r) (at ?i ?w) (not (loaded ?r ?i))))
200
201 (:action move :parameters (?r ?w1 ?w2)
202 :precondition (and (rover ?r) (wp ?w1) (wp ?w2) (at ?r ?w1))
203 :effect (and (at ?r ?w2) (not (at ?r ?w1))))
```

204 Given an AI planning problem  $(P, A, s_0, obj)$  in the logical representation,  
 205 the *state space* of the problem materializes as a transition system  $\Sigma = (S, A, c, \gamma)$   
 206 where  $S \subseteq 2^P$  is a subset of states on  $P$ . Then, solving an AI planning problem  
 207 means finding a plan to go from the initial state to a state which satisfies the  
 208 objective.

209 **Definition 4** (Solution and Optimal Solution). *A solution to an AI planning*  
 210 *problem  $(P, A, s_0, obj)$  is a plan  $\pi$  such that  $obj \subseteq \gamma(s_0, \pi)$ . The cost of a*  
 211 *solution is simply  $c(\pi)$ . An optimal solution of  $(P, A, s_0, obj)$  is a plan  $\pi^*$  such*

212 that  $\pi^* = \arg \min c(\pi)$ .

213 Note that there can be multiple optimal solutions. Finding solutions to an  
214 AI planning problem is known to be NP-hard [20] and most classical AI planners  
215 combine a search engine based on variations of the A\* algorithm with heuristics  
216 and forward/backward strategies to explore the state space.

217 The following code gives an example of a typical AI planning problem where  
218 an initial state is described with one free rover **r1**, one item **i1** and two waypoints  
219 **w1,w2**, where both **r1** and **i1** are at **w1** and the objective is to find an optimal  
220 plan which leads the system to a state where the item **i1** is transported to **w2**  
221 and the rover **r1** is at **w1**.

222 **Example 4.** (*A planning problem*)

```
223 (define (problem rover1) (:domain rover) (:objects r i w1 w2)
224 (:init (rover r) (free r) (item i) (wp w1) (wp w2)
225 (at r w1) (at i w1))
226 (:goal (and (at i w2) (at r w1))))
```

227 Using a typical AI planner, for example, the online planner **planning.domains**<sup>2</sup>  
228 [34], will give the following optimal solution:

229 **Example 5.** (*An optimal solution*)

```
230 [(load r w1 i), (move r w1 w2), (unload r w2 i), (move r w2 w1)]
```

## 231 2.2. Metamorphic Testing (MT)

232 MT aims at (i) checking the expected outputs of a system under test by using  
233 user-defined properties of the system, called *metamorphic relations* and (ii)  
234 generating *follow-up test cases* by using these relations [8, 9]. In the following,  
235 we formalize MT by reusing the notations of [47].

236 **Definition 5** (Metamorphic Relation (MR)). *Let  $f$  be a target function, an MR*  
237 *is a necessary property over a sequence of multiple inputs  $\langle x_1, \dots, x_n \rangle$  ( $n \geq 2$ )*

---

<sup>2</sup>Online at: <http://planning.domains/>

238 and their corresponding outputs  $\langle f(x_1), \dots, f(x_n) \rangle$ . MR is expressed as a relation  
 239  $\mathcal{R} \subseteq X^n \times Y^n$ .

240 **Definition 6** (Follow-up test cases). *Let the following*  
 241  $\mathcal{R}(x_1, \dots, x_n, f(x_1), \dots, f(x_n))$  *be an MR, then the sequence of inputs and their*  
 242 *corresponding outputs defines the set of source inputs. For each source in-*  
 243 *put  $\mathcal{S}(x) \in \mathcal{R}$ , a follow-up test case  $\mathcal{F}$  is derived by applying a possibly non-*  
 244 *deterministic transformation function  $T$  to the input of the source test case  $x$ :*  
 245  $\mathcal{F}(x) = f(T(x))$ . *The transformation function  $T$  is constructed such that the*  
 246 *follow-up test case  $\mathcal{F}$  fulfills the necessary property of  $\mathcal{R}$ .*

247 Let  $P$  be an implementation of  $f$  and  $\mathcal{R}$  be the MR  $\mathcal{R}(x_1, \dots, x_n, f(x_1), \dots, f(x_n))$ ,  
 248 then applying MT for  $P$  involves the following steps:

- 249 1. Define  $\mathcal{R}'$  by replacing  $f$  by  $P$  in  $\mathcal{R}$ .
- 250 2. Given a sequence of source test cases  $\langle x_1, \dots, x_k \rangle$ , execute them to obtain  
 251 their respective outputs  $\langle P(x_1), \dots, P(x_k) \rangle$ . Construct and execute a se-  
 252 quence of follow-up test cases  $\langle x_{k+1}, \dots, x_n \rangle$  according to  $\mathcal{R}'$  and obtain  
 253 their respective outputs  $\langle P(x_{k+1}), \dots, P(x_n) \rangle$ .
- 254 3. Examine the results with reference to  $\mathcal{R}'$ . If  $\mathcal{R}'$  is not satisfied, then this  
 255 MR has revealed a fault in the implementation of  $P$ .

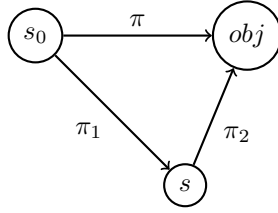
256 Using  $\mathcal{R}$  as MR, the follow-up test case  $\mathcal{F}$  with test input  $T(x)$  can be  
 257 generated using so-called *generators*. As indicated in Def.6, the transformation  
 258  $T$  is not necessarily a deterministic function, and it can lead to the production  
 259 of multiple test cases.

### 260 3. Mutation-guided MT of AI Planners

261 This section presents an automated mutation-guided MT method dedicated  
 262 to the detection of non-optimal AI planners.

263 *3.1. An MR for Detecting Non-Optimality*

264 We base our MR on a well-known property of optimal plan which establishes  
 265 that, for a given supposedly optimal initial plan, there cannot be strictly less  
 266 costly plans obtained by transiting by another reachable state than the initial  
 267 one. As illustrated below, if  $\pi$  is the initial plan returned by the planner under  
 268 test (PUT), then the cost of any plan  $\pi_1$  transiting by a state  $s$  added to the  
 269 cost of any plan  $\pi_2$  cannot be less than the cost of  $\pi$ .



270  
 271 For a given implementation *planner*, an MR for testing non-optimality can be  
 272 formally defined as follows:

**Definition 7** (*MR over planner w.r.t. s*). Let  $p = (P, A, s_0, obj)$  be an AI planning problem,  $s$  be a reachable state ( $\neq s_0$ ) and *planner* be a supposedly optimal AI planner, then an MR for  $p$  w.r.t.  $s$  is

$$\langle p, p' \rangle \implies c(\text{planner}(p)) - c(\text{planner}(p')) \leq c([a_1, \dots, a_k])$$

where  $[a_1, \dots, a_k]$  is a plan from  $s_0$  to  $s$ ,  $p'$  is defined as the AI planning problem  $p' = (P, A, s, obj)$  and  $\text{planner}(p)$  represents the call of *planner* on the problem  $p$ . If we consider only unitary cost (i.e., the cost of each action equals 1.) then the formula simplifies as

$$\langle p, p' \rangle \implies (|\pi| - |\pi'|) \leq k$$

273 where  $k$  is the number of actions performed to reach  $s$  from  $s_0$ ,  $\pi$  (resp.  $\pi'$ ) is  
 274 a plan generated by *planner* from  $s_0$  (resp.  $s$ ).

275 Note that any reachable state (but the initial state of the problem) is eligible  
 276 to reveal the non-optimality of the planner. Other MRs can be used for testing

277 AI planners, such as, for instance, relations that change the order of defined  
278 actions in the domain definition file or the order of objects in the problem defi-  
279 nition file, but these relations would be useless to test for optimality as most AI  
280 planners just preprocess their inputs to put actions and objects in a predefined  
281 similar order. The MR from Def. 7 captures the essence of optimality issues  
282 in AI planning. Two main approaches can be used to explore the state space:  
283 progressing in forward chaining from the initial state (it corresponds to reaching  
284  $s$  and asking the PUT to produce  $\pi_2$ ), or regressing in backward chaining from  
285 the objective (thus reaching  $s$  from  $obj$  and asking the PUT to produce  $\pi_1$ ).  
286 Note that we can also choose  $s$  arbitrarily and let the PUT generate both  $\pi_1$   
287 and  $\pi_2$  but this approach would be highly inefficient as the likelihood to choose  
288 a reachable state is very low. In our work, we use only the former approach, as  
289 it explores reachable states only. Indeed, the backward approach that was also  
290 implemented and evaluated leads to inefficiently exploring too many unreach-  
291 able states. We then gave up on this backward exploration approach. Note  
292 that the absence of violation of this MR does not ensure that the AI planner  
293 under test is always optimal as 1) the definition is related to selected planning  
294 tasks which means that other tasks may reveal undetected optimality faults; 2)  
295 checking the MR for all possible states  $s$  is intractable as exploring the whole  
296 state space is usually impossible and only a selection of states can be tested.  
297 Furthermore, another limitation of this approach lies in the inconsistency of the  
298 subsequent follow-up planning tasks. For some states  $s$ , the plan  $\pi_2$  does not  
299 necessarily exist and then checking the MR for  $s$  is useless.

### 300 3.2. Mutation Adequacy Test Case Selection

301 Generating follow-up test cases with the proposed MR (Def. 7) involves 1)  
302 generating states by exploring the state space from the initial state (using the  
303 forward chaining approach); 2) selecting the states which are the most error-  
304 prone; and 3) launching of new planning tasks from these states. By observing  
305 that there is no obvious relation between states and optimality fault detection,  
306 we then face a *search problem* without benefiting from any guidance on how to

---

**Algorithm 1** Mutation adequacy scoring computation

---

**Input:**  $MUT$ : a set of mutated AI planners,  $Task$ : an AI planning problem**Output:** Scores: List of scores

```
1:  $Costs \leftarrow []$ 
2: for all  $m \in MUT$  do ▷ Collect source costs
3:    $\pi \leftarrow execute(m, Task)$ ;
4:    $Costs.append(c(\pi))$ 
5: end for
6:  $Nodes \leftarrow generateStates(Task)$ 
7: for all  $(s, c) \in Nodes$  do ▷ Mutation score per state
8:    $score \leftarrow 0$ ;
9:    $Subtask \leftarrow PlanningTask(Task, s)$ 
10:  for all  $m \in MUT$  do
11:     $\pi \leftarrow execute(m, Subtask)$ 
12:    if  $Costs[m] > c + c(\pi)$  then  $score \leftarrow score + 1$ 
13:    end if
14:  end for
15:   $Scores.append(score)$ 
16: end for
17: return  $Scores$ 
```

---

307 select states. Our approach to tackling this problem is to guide state selection  
308 by using mutation testing. We guide the state selection by scoring the states  
309 with their ability to kill mutated AI planners. A mutant is killed by a state if  
310 the subsequent test case violates the MR.

311 Algorithm 1 details our mutation adequacy scoring computation. The al-  
312 gorithm starts by executing all AI mutated planners on the input AI planning  
313 problem and saves the costs of their generated plans, i.e., the *source costs* (lines  
314 1-5). Then, the state space of the problem is explored to collect potential  
315 states useful for generating follow-up test cases (line 6). Each returned state  
316  $s$  is labelled with the cost  $c$  needed to reach it from  $s_0$ . Next, we iterate over

317 the collected pairs  $(s, c)$  to compute their mutation adequacy score by first in-  
318 stantiating the associated follow-up AI planning problem (line 7-9), second, by  
319 executing all mutated AI planners (line 10-11) on this newly generated test  
320 problem, and finally by retrieving the costs of their generated plans. Then, the  
321 algorithm checks the MR and counts the number of its violations, leading to the  
322 adequacy score (line 12).

323 The analysis of Alg. 1 shows that its complexity is  $O(|MUT| \cdot |Exec| \cdot |Nodes|)$   
324 where  $|MUT|$  denotes the number of mutated AI planners,  $|Exec|$  denotes the  
325 number of calls to an AI planner and  $|Nodes|$  denotes the number of generated  
326 states. Note that all mutated AI planners are invoked once for solving each  
327 generated AI planning problem (for each newly generated state). Hence, in the  
328 proposed framework, it is necessary to limit the number of generated states  
329 ( $|Nodes|$ ) such that the complexity of the computation remains bounded. Note  
330 also that even though the cost of the state generation (call to *generateStates* in  
331 line 6) can vary depending on the chosen exploration strategy and the considered  
332 AI planning problem, the overall time complexity of Alg. 1 is dominated by the  
333 calls to the mutated AI planners ( $|Exec|$ ). Luckily, these calls can be artificially  
334 limited by setting up acceptable time-outs.

### 335 3.3. Workflow of the MT framework

336 The overall workflow of our mutation-guided MT method is shown in Fig-  
337 ure 3.3. In addition to the AI planner under test PUT, the method takes as  
338 inputs *source input* (the initial planning problem),  $N_{max}$  (the maximum number  
339 of follow-up test cases) and a set of mutated AI planners *mutants*.

340 **Initialization** The PUT is executed on *source input* and the cost of the gener-  
341 ated plan, *source output*, is computed.

342 **State Generation** This step explores the state space and gathers visited states.  
343 As said in Section 3.1, the process starts from the initial state and progresses  
344 through the state space by using forward chaining. By default, the exploration  
345 performs a breadth-first search, since it discovers the states with the shortest  
346 paths. The states are then saved and labelled with the cost of the plan used to

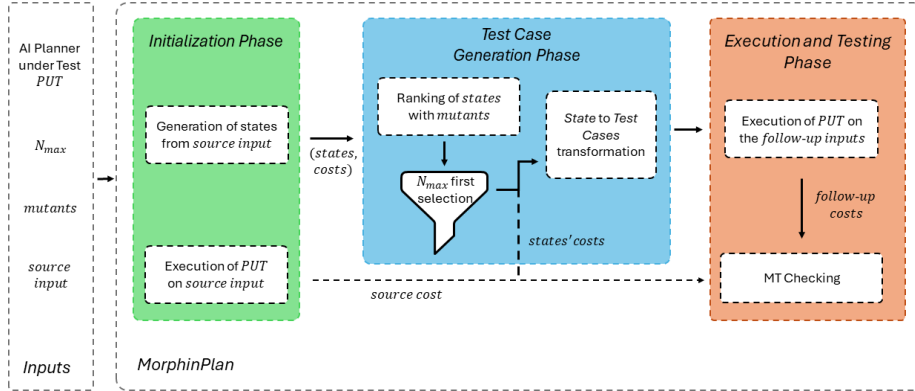


Figure 2: Overview of the MT Framework. As inputs, it receives the AI Planner under test  $PUT$ , the maximum number of follow-up test cases  $N_{max}$ , a set of mutated AI planners  $mutants$  and an initial planning task  $source\ input$  (expressed in PDDL). The boxes summarise the main stages that are sequentially executed: initialization and state generation (in green), state selection and their transformation to follow-up test cases (in blue), MT testing of  $PUT$  (in orange). Bold arrows denote the workflow, while dashed arrows note computed data that is reused in later phases.

347 reach them. Note that, as the state space can be very large for some problems,  
 348 the generation must be limited. In our implementation, we stop the generation  
 349 when a user-specified number of states has been reached, but other stopping  
 350 criteria can be used, such as a pre-selected time-out.

351 **State Selection** The generated states are first scored and ranked *w.r.t.* their  
 352 respective adequacy score on the set of mutated AI planners. Then, the method  
 353 returns the  $N_{max}$  first states whose costs are lower than the  $source\ output$ 's one.

354 **Follow-up Input Generation** Each selected state is used to create a new  
 355 planning problem, i.e., a follow-up test case, by producing a new PDDL problem  
 356 file. Given a generated state  $s$  and the  $source\ input$   $(P, A, s_0, obj)$ , the follow-up  
 357 test case is the planning problem  $(P, A, s, obj)$ .

358 **Execution** This step eventually executes the  $PUT$  on the set of follow-up test  
 359 cases, retrieves its outputs, called *follow-up outputs*, and returns their respective  
 360 costs.

361 **Validation** The MR is checked for every *follow-up output* cost, that is, its sum  
 362 with the cost of the related state must be higher or equal than the cost of the

363 *source output.*

## 364 **4. Implementation and Experimental Setup**

### 365 *4.1. Implementation*

366 We implemented our MT framework for detecting optimality faults in AI  
367 planners in a tool called MORPHINPLAN. The implementation avoids redundant  
368 computations as much as possible by saving all the generated states for each  
369 planning task as well as the results of the mutated AI planners. The tool  
370 is implemented in SICStus Prolog [7] in approximately 4KLOC and it can test  
371 any PDDL2.1-based AI planner, provided that there are no durative actions and  
372 preconditions, and no continuous effects. In MORPHINPLAN, mutation testing  
373 is used at two distinct levels: i) to implement the mutation-based test case  
374 selection strategy and ii) to test MORPHINPLAN w.r.t. its capability to detect  
375 optimality faults in generated plans by AI planners. Therefore, for both cases,  
376 we created 22 non-optimal by-design AI planners, i.e., *mutated planners*, from  
377 a *reference AI planning implementation* to create source-mutated AI planners.  
378 These mutants are precisely described in Section 4.3 as designing functional  
379 mutants which exhibit optimality faults is not a trivial task.

380 To implement our mutation adequacy test case selection strategy, we con-  
381 sidered a state generation based on a breadth-first search exploration with a  
382 maximum boundary of 500 states. In our experiments, we observed that vary-  
383 ing this number by 20% has no impact on the results.

384 Following the Open Science Policy, we provide the source code of MORPHIN-  
385 PLAN as well as all the needed material to replicate the results presented in the  
386 paper. The MORPHINPLAN repository is [https://github.com/QuentinMaz/](https://github.com/QuentinMaz/MTFramework/tree/feature/fse)  
387 [MTFramework/tree/feature/fse](https://github.com/QuentinMaz/MTFramework/tree/feature/fse) and the 22 mutated planners are available in  
388 <https://github.com/QuentinMaz/MutatedAIPlanner/tree/develop>.

### 389 *4.2. Experimental Setup*

390 Our experiments consist of executing MORPHINPLAN on AI planners under  
391 test (PUTs) on a selection of PDDL AI planning problems. To address as best

392 possible the huge variety of AI planning problems, we selected a subset of PDDL  
393 problems from those used in the international planning competition (IPC).<sup>3</sup> The  
394 problems were selected on 1) their diversity by picking up problems coming from  
395 different application areas (i.e., domains such as block-world, robotics, logistics,  
396 etc.) and 2) the current restrictions of our reference AI planner implementation  
397 which supports only PDDL2.1 with no durative and continuous actions. The  
398 selection resulted in 19 problems from 10 distinct domains, which are listed in  
399 Table 1 with their reference.

#### 400 4.3. Mutated AI Planners

401 To support mutation-guided MT and to perform mutation testing of MOR-  
402 PHINPLAN, we created hard-coded non-optimal-by-design AI planners from Sasak’s  
403 PPL, an existing open-source Prolog-based AI Planning library [42]. Since  
404 applying standard mutation operators (e.g., Relational or Boolean Operator  
405 Replacement) to AI planner implementations had almost invariably led us to  
406 unusable planners (i.e., planners that cannot generate any plan), we manually  
407 designed 22 mutated AI planners intending to generate non-optimal plans in  
408 some cases. Note that creating such mutated AI planners is not trivial, as it  
409 requires the introduction of subtle faults that lead the planners to produce a  
410 valid but not optimal plan, i.e., a plan which satisfies the objective but has a  
411 strictly higher cost than an optimal plan. Our idea was to implement a forward  
412  $A^*$  search that does not reopen the closed set (i.e., states that have already  
413 been visited), and where the  $f$  function (used to prioritize the open set) and  
414 the  $h$  function (used to approximate the distance of a state to the goal) are  
415 mutated. We developed 5  $f$  and 6  $h$  functions, which are presented in Table 2.  
416 We then compared the plans of the 30 ( $f$ ,  $h$ ) combinations to filter *equivalent*  
417 *mutants* [31]. In our context, a mutated planner is an *equivalent mutant* [31] if  
418 and only if, despite the introduced modification, the mutated planner generates  
419 only optimal plans for all the considered planning problems, as done by the

---

<sup>3</sup><https://www.icaps-conference.org/competitions/>

Domain	Description
airport	Ground traffic control on airports [24].
blocks	Blocks world: Stackable blocks must be re-assembled on a table [3].
gripper	A robot with two grippers is tasked to take balls from one room to another [33].
elevator	Elevator control: Transport passengers with an elevator under side constraints [3].
openstacks	Minimum maximum simultaneous open stacks: Order production tasks such that there is only a minimal number of open orders [16].
pegsol	Peg Solitaire game [35].
psr-small	Power Supply Restoration: Resupply number of lines in a faulty electricity network [24].
tpp	Travelling Purchaser Problem [16].
transport	Logistics planning problem used in the IPC 2008 competition [35].
travel	Travel logistics involving side constraints.

Table 1: Overview of the planning domains used in the experimental evaluation.

420 original implementation. After filtering, we ended up with 22 distinct mutated  
 421 planners, as shown in Table 3.

#### 422 4.4. Baseline Strategies

423 We have identified three simple strategies, alternative to our mutation-based  
 424 test case selection method, as comparison baselines. These strategies replace the  
 425 state generation and/or the state selection components. For selecting  $N$  states  
 426 as follow-up test cases, we consider the following strategies:

427 `bfs_det` Breadth-first search from  $s_0$  and selection of the  $N$  first-found states;

428 bfs\_ran Breadth-first search from  $s_0$  and random selection of  $N$  visited states  
429 among 500 generated states;

430 wal\_ran Iterative random walk from  $s_0$  repeated  $N$  times. The maximum length  
431 of each walk is limited to the cost of the *source plan*;

432 These baselines are complementary as they correspond to three non-costly  
433 strategies for generating/selecting states in the state space. The first one sys-  
434 tematically explores the state space from  $s_0$  and evaluates whether further ex-  
435 ploration is needed. The second one challenges whether all states are equally  
436 capable of revealing optimality faults. Eventually, the third baseline strategy  
437 evaluates the general relevance of a dedicated state generation and selection  
438 functionality, while a simple random walk exploration can be used.

## 439 5. Experimental Evaluation

### 440 5.1. Research Questions

441 The evaluation of our mutation-guided MT framework led us to answer the  
442 three following research questions (RQs):

443 **RQ1** How effective is the proposed framework at detecting optimality faults in  
444 mutated AI planners?

445 **RQ2** How efficient is *mutation adequacy test selection* at detecting optimality  
446 faults, as compared to identified baseline strategies?

447 **RQ3** Can the framework detect optimality faults in off-the-shelf AI planners?

448 We conducted two main experiments to answer the three RQs. For the first  
449 experiment, which aims to answer RQ1 and RQ2, we randomly split the set of  
450 mutated AI planners: some of them become PUTs to evaluate the capability of  
451 MORPHINPLAN to reveal optimality faults in AI planners, and the others serve  
452 to guide the mutation adequacy test selection strategy. To mitigate bias on the  
453 mutated AI planners' selection for one or the other task, we report averaged

454 results over 20 executions with different random splits on the set of mutated  
455 planners.

456 Regarding the second experiment, which answers RQ3, we used MORPHIN-  
457 PLAN to test existing off-the-shelf planners and used all the mutated AI plan-  
458 ners available to guide the test selection. Since MORPHINPLAN can test any  
459 PDDL2.1-based AI planner, we decided to source the PUTs from two versions  
460 of a reference AI planner, namely the PPL planner described in Section 4.3  
461 and several distinct configurations of a well-recognized implementation, namely  
462 Fast Downward<sup>4</sup> [23], which won one of the IPC 2018 competition track. Fast  
463 Downward is one the most used planning systems in research, with more than  
464 a decade of contributions from the AI planning community.<sup>5</sup>

465 For both experiments, we ran MORPHINPLAN with  $N_{max} = 4$  which means  
466 that only 4 follow-up test cases were considered for detecting non-optimal plans.  
467 In subsequent runs of MORPHINPLAN dedicated to analysing the importance of  
468 the  $N_{max}$  value in the process, we varied that value from 1 to 60 (see Figure 4).  
469 We systematically measured the runtime needed for the generation and selection  
470 of test cases in all cases. We excluded the time taken by the mutated AI planners  
471 as 1) it is artificially limited by a parameterized time-out value and 2) to keep  
472 the comparison fair w.r.t. the other strategies which do not call the mutated  
473 AI planner.

474 All our experiments were performed on a Linux machine (Ubuntu 22.04.3  
475 LTS) equipped with an AMD Ryzen 9 3950X 16-Core processor and 32GB of  
476 RAM.

## 477 5.2. RQ1: Effective detection of mutated AI planners

478 To answer RQ1, we measure mutation scores and running times achieved by  
479 MORPHINPLAN on the aforementioned 19 AI planning problems. The mutation  
480 scores are reported as percentages of mutants detected as non-optimal by at

---

<sup>4</sup>Fast Downward: <https://www.fast-downward.org>

<sup>5</sup>Brief history of Fast Downward: [https://github.com/aibasel/downward/tree/main?  
tab=readme-ov-file#contributors](https://github.com/aibasel/downward/tree/main?tab=readme-ov-file#contributors).

481 least one follow-up test case. Table 4 1<sup>st</sup> column presents the average results  
 482 over the 20 subsets of validation/selection mutants obtained with MORPHIN-  
 483 PLAN, as well as the execution time. We observe that MORPHINPLAN gives  
 484 good averaged mutation scores on most problems. More precisely, 15 over 19  
 485 problems have revealed non-optimality in the AI mutated planners, with aver-  
 486 aged mutation scores greater than 80% with acceptable deviation (around 5%  
 487 in most cases). gripper01, psr-small08 and tpp03 have mutation scores close to  
 488 the threshold of 80% (around 76%), with deviations around 9%. openstacks01  
 489 has only a mutation score of circa 27%, which will be discussed in RQ2, along  
 490 with the results of the baseline strategies. As the mutated AI planners used  
 491 to guide MORPHINPLAN have been randomly selected among the 22 available  
 492 mutants, our experiment shows that the mutation-guided MT framework is  
 493 agnostic w.r.t. the selection of mutants. It also shows that the framework is  
 494 efficient in detecting non-optimality in mutated AI planners by using any PDDL  
 495 domain/problem. Only one problem (i.e., openstacks01) exhibits a low muta-  
 496 tion score, while almost all others achieve scores greater than 80%. Regarding  
 497 execution time, Table 4 shows that the execution time of MORPHINPLAN re-  
 498 mains acceptable as it ranges between 20s and 214s with an average of 60s.  
 499 This mean value is an acceptable execution time for an auto-test procedure to  
 500 be embedded into an autonomous system. Note that this mean value can also be  
 501 considered regarding the time needed to execute the planner to let it compute  
 502 an optimal plan. Even if the time needed is artificially bounded, it is usually  
 503 set to longer time slots (typically 360s or more).

504 *5.3. RQ2: Efficiency of mutation adequacy test selection as compared to simpler*  
 505 *strategies*

506 In Table 4, we report the average mutation scores achieved by the three  
 507 baseline strategies, namely bfs\_det, bfs\_ran and wal\_ran along with their execu-  
 508 tion times. These results show that MORPHINPLAN performs best on all of the  
 509 19 problems, with a mean of 84.7%. Looking at the mean results, we see that  
 510 both bfs-based strategies achieve close mutation scores of  $\sim 52\%$  while wal\_ran

511 achieves only  $\sim 34\%$ . The result obtained in openstacks01 shows that all four  
 512 strategies achieve poor mutation scores. The explanation lies in the hardness of  
 513 killing some mutants with this problem. As the two non-deterministic strate-  
 514 gies (i.e., bfs\_ran and wal\_ran) outperform bfs\_det, it seems that the mutants  
 515 killed by all these three strategies are similar, and the other mutants remain  
 516 alive for this problem. To refine our analysis, we also looked into the per-  
 517 centage of generated follow-up test cases that lead to a violation of the MR  
 518 (i.e., test cases that can detect the generation of non-optimal plans). Figure 3  
 519 shows the averaged results over the 19 problems compactly. We observe that  
 520 MORPHINPLAN achieves a much higher percentage (higher than 70%) than the  
 521 three other strategies (less than 28%). It means that the ability to generate  
 522 fault-revealing test cases is much higher for MORPHINPLAN. Regarding run-  
 523 ning times, the mean time needed by the alternative strategies is much less for  
 524 bfs\_det and bfs\_ran than for MORPHINPLAN. It is understandable as MORPHIN-  
 525 PLAN explores the state space up to 500 states, while the other strategies limit  
 526 their exploration to  $N_{max}$ . However, as discussed in Section 5.2, the mean time  
 527 taken by MORPHINPLAN is acceptable as it never exceeds the auto-test time  
 528 limit of 360s.

529 By studying the influence of the  $N_{max}$  value over the mutation scores re-  
 530 sults, as shown in Figure 4 (right plot), our experiment shows that the baselines'  
 531 scores increase with the number of test cases executed. However, MORPHIN-  
 532 PLAN quickly achieves a high mutation score with only a few follow-up test  
 533 cases, and stabilizes its results from  $N_{max} \geq 10$ . This swift convergence em-  
 534 phasizes the overall efficiency observed in Figure 3 and reveals that the number  
 535 of test cases needed to achieve high mutation scores is much less than the three  
 536 competitive baselines. Regarding RQ2, not only the mutation adequacy state  
 537 selection strategy leads to better non-optimal behaviour detection compared to  
 538 the baselines, but also lowers the needed number of follow-up test cases, under-  
 539 lining the efficiency of mutation adequacy test selection for our MT framework.  
 540 However, it is stressed that the time needed to execute that strategy is much  
 541 higher than the one needed for simpler strategies. It is thus crucial to control

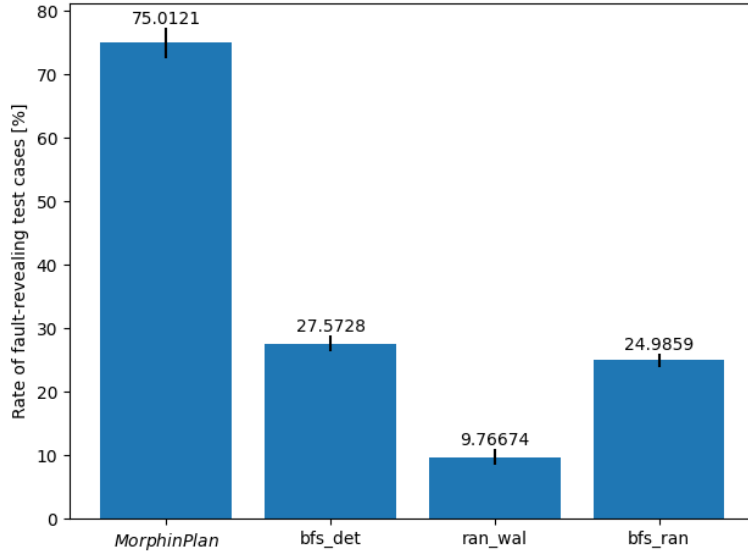


Figure 3: Percentage of MR-violating test cases for the different strategies. The results are averaged across the 19 problems considered.

542 it by appropriate time-out values.

543 *5.4. RQ3: Detection of non-optimality in off-the-shelf planners*

544 To answer RQ3, we synthesize the results obtained by our second experi-  
 545 ment, which are reported in Table 5. Regarding off-the-shelf planners, we con-  
 546 sidered three configurations of the Sasak’s PPL (V0, V1- $A_{h0}$  and V1- $A_{hmax}$ )  
 547 and a total of 17 configuration settings for Fast Downward, i.e., using differ-  
 548 ent internal parameters. Among these 17 configurations, 7 are using heuristics  
 549 and algorithms that lead the planner to generate only optimal plans (column  
 550 “Opt”) and 10 are using non-admissible heuristics. We considered these confi-  
 551 gurations as other baselines to check whether MORPHINPLAN can detect these  
 552 non-optimal planners. Interestingly, we considered V0 of Sasak’s PPL because  
 553 it contains a subtle known optimality fault that was later corrected in V1. We  
 554 wanted to see if MORPHINPLAN could retrieve that bug. Table 5 shows the  
 555 results of non-optimality detection per AI planning problem. First, note that

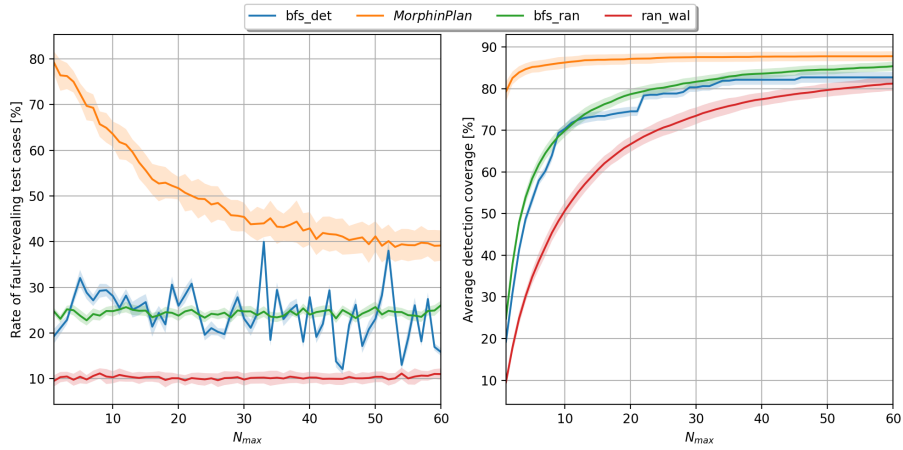


Figure 4: Evolution of the rates of fault-revealing test cases (left) and mutation scores (right) w.r.t.  $N_{max}$  values.

556 Sasak’s PPL could not handle some problems due to the absence of support for  
557 certain features of PDDL2.1 (column N/A) in that implementation. Second, the  
558 7 optimal configurations of Fast Downward have been successfully tested with  
559 MORPHINPLAN without detecting any non-optimality behaviour. Regarding the  
560 known bug of PPL V0, it was revealed only on two problems from a single do-  
561 main (elevator). It indicates that MORPHINPLAN shall be used with a sufficient  
562 variety of PDDL domain/problems as non-optimality bugs can sometimes only  
563 be revealed by a few domain definitions. Third, non-optimality behaviours were  
564 detected for all the 10 remaining configurations of Fast Downward (at least one  
565 ‘Check’ in all columns). This result is interesting, as the variety of problems,  
566 combined with the capability of MORPHINPLAN to explore the state space, leads  
567 to the systematic detection of existing non-optimal behaviours in AI planning.  
568 Regarding RQ3, the results of this second experiment show that MORPHINPLAN  
569 systematically detects non-optimal behaviours in the two selected AI planners.

## 570 6. Threats to Validity

571 As *internal threats-to-validity*, we identified the creation of mutated AI plan-  
572 ners for guiding the MT framework and the choice of the maximum number of  
573 follow-up test cases. To cope with the first threat, we randomized the selec-  
574 tion of mutants used by MORPHINPLAN and provided only average results, plus  
575 standard deviations. By doing this, our results show that MORPHINPLAN is ag-  
576 nostic to the usage of mutants to guide the test case selection. However, the 22  
577 mutated planners were manually created by using a single idea (as explained in  
578 Section 4.3) and thus they may not be representative enough of the diversity of  
579 potential faults found in AI planning implementation. To cope with this threat,  
580 our second experiment dealt with off-the-shelf planners only, without any evalu-  
581 ation of the mutated AI planners. Regarding the maximum number of follow-up  
582 test cases, the results of Figure 4 indicate that the framework is dependent on  
583  $N_{max}$  and the choice of this value may largely influence the results of any test  
584 case selection strategy. Even though more experiments are needed to further  
585 analyse this question, the preliminary results of Figure 4 are very encouraging.

586 As *external threats-to-validity*, we identified i) the selection of PDDL prob-  
587 lems to perform our experiments, ii) the selection of planners to test and, iii)  
588 the applicability of MORPHINPLAN only to optimal AI planner. As explained  
589 in Section 4.2, our selection was objective but still, other PDDL domains and  
590 problems could lead to different results as nothing guarantees that our selection  
591 is representative enough of the huge diversity of AI planning problems. Also,  
592 our second experiment included only two AI planners in 20 distinct configura-  
593 tions. Even though these planners are off-the-shelf planners, the results could  
594 be different for other planners. By using a single MR dedicated to finding non-  
595 optimality issues, our approach focuses on testing optimal AI planning systems.  
596 Of course, other realistic planners, not necessarily seeking optimality but just  
597 validity of the generated plan, remain outside the scope of MORPHINPLAN. To  
598 cope with these threats and limitations, more experiments are needed to confirm  
599 the obtained results and the usage of more MRs are needed.

## 600 7. Related Work

601 The validation and verification of AI planning systems is an important issue  
602 that has been addressed from different angles since the last century [4]. Some ap-  
603 proaches have focused on the testing of domain and problem definition, chasing  
604 faults in PDDL programs [26, 41, 22, 32]. A major approach for the verification  
605 of safety properties in domain specification is based on model-checking methods,  
606 especially for industrial application cases [37, 6, 46, 45, 1]. The general idea is  
607 to perform bounded model checking on the underlying state space associated  
608 with the specified PDDL domain. Interestingly, other approaches have formu-  
609 lated test case generation as AI planning problems [25, 5]. Our approach differs  
610 largely from these works in that it focuses on the testing of the AI planners for  
611 detecting non-optimal generated plans. Under the hypothesis that domain and  
612 problem specification are verified by other means, our contribution lies in the  
613 usage of mutation-guided MT for detecting optimality-related faults.

614 MT has been successful in testing general and autonomous software systems  
615 [44, 30, 29, 14, 2, 15, 49]. The essential contribution of these works lies in the  
616 definition of advanced MRs for complex learning systems. Unlike our approach  
617 which focuses on using one MR for checking the optimality of an AI planner,  
618 these works contribute to more general system testing without questioning the  
619 generation of follow-up test cases for AI planners. Recently, MT [18] was used  
620 to test learned action policy used in games and planning. The idea was to  
621 exploit state relaxation, an easy abstraction mechanism used in AI planning, to  
622 construct MRs.

623 The closest works regarding the usage of MT for checking optimality faults  
624 in AI planning systems are by Zhang et al. [54] and by Cheng et al. [12]. In  
625 [54], the authors propose to use the distance-related triangle inequality as an  
626 MR for testing graph-based path planning algorithms. Our MR defined based  
627 on cost-driven minimisation has similarities with this triangle inequality, but it  
628 differs from it by avoiding the necessary distance definition which is necessary  
629 for path planning. Also, the approach generalizes the usage of MT by working

630 directly at the level of PDDL for any domain definition (not only for path plan-  
631 ning) which allows us to deal with a much larger class of problems. Note that  
632 MORPHINPLAN enables a dynamic exploration of the underlying search space  
633 without requiring its explicit construction. In [12], the goal is to test learnt  
634 policies for automated driving applications by using MT. The employed MR is  
635 based on mutating an input situation in such a way that optimality shall be  
636 preserved. Thus, it resembles our MR, which looks for derived states where op-  
637 timality can easily be compared with the one provided by an initial planner run.  
638 However, our approach differs in its goal (testing a PDDL-based AI planning  
639 system) and approach (using mutated versions of a reference implementation).  
640 Unlike [54, 2, 55, 12], our approach introduces state selection strategies, includ-  
641 ing mutation adequacy selection, for generating the follow-up test cases. Despite  
642 the importance of recent findings in mutation testing [36, 48], to the best of our  
643 knowledge, it is the first attempt to use mutation testing for verifying AI plan-  
644 ning systems. The difficulty of creating mutants for complex and optimized AI  
645 planners that exhibit optimality-related faults may explain the absence of such  
646 a tool in the literature.

## 647 **8. Conclusion and Future Work**

648 In this article, we have introduced a MT framework for testing optimality  
649 in AI planners with a mutation adequacy test selection strategy. Our tool  
650 MORPHINPLAN considers any PDDL planning domain and problem (written in a  
651 restricted fragment of PDDL2.1) and generates, based on the PUT's initial plan,  
652 new problems as follow-up test cases by exploring the state space of the problem.  
653 By comparing the actual results for these new problems to the cost of the initial  
654 plan, non-optimal AI planners can be detected using MT. In our framework,  
655 we score all the visited states with their ability to kill mutated AI planners,  
656 that have been previously created from a reference implementation, and we  
657 select the states having the greatest scores, i.e., the states most likely to reveal  
658 faults, to generate follow-up test cases. Our experimental evaluation shows that

659 MORPHINPLAN is effective at detecting non-optimal AI planners among both  
660 artificially created mutated AI planners and off-the-shelf, configurable planners.  
661 A key observation of our results is that mutation adequacy test case selection  
662 is well-tuned for detecting non-optimality and that it outperforms traditional  
663 breadth-first search and random strategies, provided that sufficient time can be  
664 allocated to the test generation and selection steps.

665 Our future work encompasses two directions. The first direction follows the  
666 recent research efforts in applying MT to learnt policies on Markov Decision  
667 Processes. For these models, only probabilistic AI planning approaches are rel-  
668 evant, as these processes are inherently non-deterministic. The second direction  
669 is the development of MT for AI planners based on PDDL3. In particular, we  
670 want to investigate MRs over the PDDL3 grammar (including domain definition  
671 and problem instances) for automating testing of PDDL planners.

## 672 References

- 673 [1] Mohammad Abdulaziz and Friedrich Kurz. 2023. Formally Verified SAT-  
674 Based AI Planning. In *Proc. of the AAAI Conference on Artificial Intelli-*  
675 *gence*. 14665–14673. <https://doi.org/10.1609/aaai.v37i12.26714>
- 676 [2] Jon Ayerdi, Pablo Valle, Sergio Segura, Aitor Arrieta, Goiuria Sagardui,  
677 and Maite Arratibel. 2023. Performance-Driven Metamorphic Testing of  
678 Cyber-Physical Systems. *IEEE Transactions on Reliability* 72, 2 (June  
679 2023), 827–845. <https://doi.org/10.1109/TR.2022.3193070>
- 680 [3] Fahiem Bacchus. 2001. AIPS 2000 planning competition: The fifth interna-  
681 tional conference on artificial intelligence planning and scheduling systems.  
682 *Ai magazine* 22, 3 (2001), 47–47.
- 683 [4] S. Bensalem, K. Havelund, and A. Orlandini. 2014. Verification and vali-  
684 dation meet planning and scheduling. *Int. Journal of Soft. Tools and Tech-*  
685 *nology Transfer* 16 (2014), 1–12. [https://doi.org/doi.org/10.1007/](https://doi.org/doi.org/10.1007/s10009-013-0294-x)  
686 [s10009-013-0294-x](https://doi.org/doi.org/10.1007/s10009-013-0294-x)

- 687 [5] Josip Božić, Oliver Tazl, and Franz Wotawa. 2019. Chatbot Testing Using  
688 AI Planning. In *IEEE International Conference On Artificial Intelligence*  
689 *Testing, AITest 2019*. <https://doi.org/10.1109/AITest.2019.00-10>
- 690 [6] Guillaume P. Brat, Doron Drusinsky, Dimitra Giannakopoulou, Allen  
691 Goldberg, Klaus Havelund, Michael R. Lowry, Corina S. Pasareanu, Ar-  
692 naud Venet, Willem Visser, and Richard Washington. 2004. Experimental  
693 Evaluation of Verification and Validation Tools on Martian Rover Soft-  
694 ware. *Formal Methods Syst. Des.* 25, 2-3 (2004), 167–198. <https://doi.org/10.1023/B:FORM.0000040027.28662.a4>
- 696 [7] Mats Carlsson and Et al. 2019. *SICStus Prolog User’s Manual, Release*  
697 *4.5.1*.
- 698 [8] T.Y. Chen, S.C. Cheung, and S.M. Yiu. 1998. *Metamorphic Testing: A New*  
699 *Approach for Generating Next Test Cases*. Technical Report HKUST-CS98-  
700 01. Department of Computer Science, Hong Kong University of Science and  
701 Technology, Hong Kong.
- 702 [9] Tsong Yueh Chen, Fei-Ching Kuo, Huai Liu, Pak-Lok Poon, Dave Towey,  
703 T. H. Tse, and Zhi Quan Zhou. 2018. Metamorphic Testing: A Review  
704 of Challenges and Opportunities. *Comput. Surveys* 51, 1 (2018), 1–27.  
705 <https://doi.org/10.1145/3143561>
- 706 [10] Tsong Yueh Chen and T. H. Tse. 2021. New Visions on Metamorphic  
707 Testing after a Quarter of a Century of Inception. In *Proc. of the 29th ACM*  
708 *Joint Meet. on European Soft. Eng. Conf. and Symp. on the Foundations*  
709 *of Soft. Eng. (ESEC/FSE), Aug. 23-28*. 1487–1490.
- 710 [11] Yuting Chen, Ting Su, Chengnian Sun, Zhendong Su, and Jianjun Zhao.  
711 2016. Coverage-Directed Differential Testing of JVM Implementations. In  
712 *Proceedings of the 37th ACM SIGPLAN Conference on Programming Lan-*  
713 *guage Design and Implementation (PLDI’16)*. 85–99.

- 714 [12] Mingfei Cheng, Yuan Zhou, Xiaofei Xie, Junjie Wang, Guozhu Meng, and  
715 Kairui Yang. 2024. Evaluating Decision Optimality of Autonomous Driving  
716 via Metamorphic Testing. arXiv:2402.18393 [cs.AI] [https://arxiv.org/  
717 abs/2402.18393](https://arxiv.org/abs/2402.18393)
- 718 [13] Stevo Alves de Andrade, Fatima L. S. Nunes, and Marcio Eduardo Dela-  
719 maro. 2023. Exploiting deep reinforcement learning and metamorphic test-  
720 ing to automatically test virtual reality applications. *Software Testing, Ver-  
721 ification and Reliability* (2023). <https://doi.org/10.1002/stvr.1863>  
722 Published online Sep. 23rd.
- 723 [14] Yao Deng, Guannan Lou, Xi Zheng, Tianyi Zhang, Miryung Kim, Huai  
724 Liu, Chen Wang, and Tsong Yueh Chen. 2021. BMT: Behavior Driven  
725 Development-based Metamorphic Testing for Autonomous Driving Models.  
726 In *2021 IEEE/ACM 6th International Workshop on Metamorphic Testing  
727 (MET)*. 32–36. <https://doi.org/10.1109/MET52542.2021.00012>
- 728 [15] Yao Deng, Xi Zheng, Tianyi Zhang, Huai Liu, Guannan Lou, Miryung  
729 Kim, and Tsong Yueh Chen. 2022. A Declarative Metamorphic Testing  
730 Framework for Autonomous Driving. *IEEE Transactions on Software En-  
731 gineering* (2022). <https://doi.org/10.1109/TSE.2022.3206427>
- 732 [16] Yannis Dimopoulos, Alfonso Gerevini, Patrik Haslum, and Alessandro  
733 Saetti. 2006. The benchmark domains of the deterministic part of IPC-  
734 5. *Abstract Booklet of the competing planners of ICAPS-06* (2006), 14–19.
- 735 [17] Yefim Dinitz and Shay Solomon. 2007. On Optimal Solutions for the Bottle-  
736 neck Tower of Hanoi Problem. In *Int. Conf. on Current Trends in Theory  
737 and Practice of Computer Science (SOFSEM) (LNCS, Vol. 4362)*. 248–  
738 259.
- 739 [18] Hasan Ferit Eniser, Timo P. Gros, Valentin Wüstholtz, Jörg Hoffmann,  
740 and Maria Christakis. 2022. Metamorphic Relations via Relaxations: An  
741 Approach to Obtain Oracles for Action-Policy Testing. In *Proc. of the 31st*

- 742 *ACM SIGSOFT Int. Symp. on Software Testing and Analysis (ISSTA)*.  
743 <https://doi.org/10.1145/3533767.3534392>
- 744 [19] Giovani Farias, Ramon Fraga Pereira, Lucas Hilgert, Felipe Meneguzzi,  
745 Renata Vierira, and Rafael Bordini. 2017. Predicting Plan Failure by Mon-  
746 itoring Action Sequences and Duration. *ADCAIJ: Advances in Distributed*  
747 *Computing and Artificial Intelligence Journal* 6, 2 (2017).
- 748 [20] Malik Ghallab, Dana Nau, and Paolo Traverso. 2016. *Automated Planning*  
749 *and Acting*. Cambridge University Press. [https://projects.laas.fr/](https://projects.laas.fr/planning/book.pdf)  
750 [planning/book.pdf](https://projects.laas.fr/planning/book.pdf)
- 751 [21] Patrik Haslum, Nir Lipovetzky, Daniele Magazzeni, and Christian Muise.  
752 2019. An Introduction to the Planning Domain Definition Language. *Syn-*  
753 *thesis Lectures on Artificial Intelligence and Machine Learning* 13 (04  
754 2019), 1–187. <https://doi.org/10.2200/S00900ED2V01Y201902AIM042>
- 755 [22] Klaus Havelund, Alex Groce, Gerard Holzmann, Rajeev Joshi, and Mar-  
756 garet Smith. 2009. Automated testing of planning models. *Lecture Notes*  
757 *in Computer Science* 5348 LNAI (2009). 5th Int. Workshop on Model  
758 Checking and Art. Intel., MoChArt 2008.
- 759 [23] Malte Helmert. 2006. The fast downward planning system. *Journal of*  
760 *Artificial Intelligence Research* 26 (2006), 191–246.
- 761 [24] Jörg Hoffmann and Stefan Edelkamp. 2005. The deterministic part of IPC-  
762 4: An overview. *Journal of Artificial Intelligence Research* 24 (2005), 519–  
763 579.
- 764 [25] Adele Howe, Anneliese Mayrhauser, and Richard Mraz. 1997. Test Case  
765 Generation as an AI Planning Problem. *Automated Software Engineering*  
766 4 (Jan. 1997), 77–106.
- 767 [26] R. Howey and D. Long. 2003. VAL’s Progress: The Automatic Valid-  
768 ation Tool for PDDL2.1 used in the International Planning Competition. In

- 769 *Proceedings of the ICAPS 2003 workshop on "The Competition: Impact,*  
770 *Organization, Evaluation, Benchmarks"*.
- 771 [27] Casidhe Hutchison, Milda Zizyte, Patrick E. Lanigan, David Guttendorf,  
772 Michael Wagner, Claire Le Goues, and Philip Koopman. 2018. Robust-  
773 ness Testing of Autonomy Software. In *2018 IEEE/ACM 40th Interna-*  
774 *tional Conference on Software Engineering: Software Engineering in Prac-*  
775 *tice Track (ICSE-SEIP), Gothenburg, Sweden*. 276–285.
- 776 [28] Upulee Kanewala, James M. Bieman, and Asa Ben-Hur. 2016. Predicting  
777 Metamorphic Relations for Testing Scientific Software: A Machine Learning  
778 Approach Using Graph Kernels. *Software Testing, Verification and Relia-*  
779 *bility* 26, 3 (May 2016), 245–269. <https://doi.org/10.1002/stvr.1594>
- 780 [29] Rui Li, Huai Liu, Guannan Lou, Xi Zheng, Xiao Liu, and Tsong Yueh Chen.  
781 2021. Metamorphic Testing on Multi-module UAV Systems. In *2021 36th*  
782 *IEEE/ACM International Conference on Automated Software Engineering*  
783 *(ASE)*. <https://doi.org/10.1109/ASE51524.2021.9678841>
- 784 [30] Mikael Lindvall, Adam Porter, Gudjon Magnusson, and Christoph Schulze.  
785 2017. Metamorphic Model-Based Testing of Autonomous Systems. *2nd In-*  
786 *ternational Workshop on Metamorphic Testing, MET 2017* (2017). [https:](https://doi.org/10.1109/MET.2017.6)  
787 [//doi.org/10.1109/MET.2017.6](https://doi.org/10.1109/MET.2017.6)
- 788 [31] Lech Madeyski, Wojciech Orzeszyna, Richard Torkar, and Mariusz Jozala.  
789 2014. Overcoming the Equivalent Mutant Problem: A Systematic Lit-  
790 erature Review and a Comparative Experiment of Second Order Mu-  
791 tation. *IEEE Transactions on Software Engineering* (2014). [https:](https://doi.org/10.1109/TSE.2013.44)  
792 [//doi.org/10.1109/TSE.2013.44](https://doi.org/10.1109/TSE.2013.44)
- 793 [32] Maurício C. Magnaguagno, Ramon Fraga Pereira, Martin D. Móre, and  
794 Felipe Meneguzzi. 2020. *Web Planner: A Tool to Develop, Visualize, and*  
795 *Test Classical Planning Domains*. In *Book Knowledge Engineering Tools*  
796 *and Techniques for AI Planning*. Springer, Chapter 11, 209–227. [https:](https://doi.org/10.1007/978-3-030-38561-3_11)  
797 [//doi.org/10.1007/978-3-030-38561-3\\_11](https://doi.org/10.1007/978-3-030-38561-3_11)

- 798 [33] Drew M McDermott. 2000. The 1998 AI planning systems competition. *AI*  
799 *magazine* 21, 2 (2000), 35–35.
- 800 [34] Christian Muise. 2016. Planning.domains. In *Proc. of Int. Conf. on Auto-*  
801 *mated Planning and Scheduling. System Demo. Track. London, UK. June*  
802 *12-17*. 242–250.
- 803 [35] IPC 2008 Organizers. 2008. IPC 2008 Website. [https://ipc08.](https://ipc08.icaps-conference.org/deterministic/index.html)  
804 [icaps-conference.org/deterministic/index.html](https://ipc08.icaps-conference.org/deterministic/index.html)
- 805 [36] Mike Papadakis, Donghwan Shin, Shin Yoo, and Doo-Hwan Bae. 2018.  
806 Are Mutation Scores Correlated with Real Fault Detection? A Large Scale  
807 Empirical Study on the Relationship between Mutants and Real Faults.  
808 In *Proc. of the 40th Int. Conf. on Software Engineering (ICSE)*. [https:](https://doi.org/10.1145/3180155.3180183)  
809 [//doi.org/10.1145/3180155.3180183](https://doi.org/10.1145/3180155.3180183)
- 810 [37] John Penix, Charles Pecheur, and Klaus Havelund. 1999. Using Model  
811 Checking to Validate AI Planner Domain Models. In *Proc. of the 23rd*  
812 *Annual Software Engineering Workshop, NASA Goddard*. [https://ntrs.](https://ntrs.nasa.gov/citations/19990054671)  
813 [nasa.gov/citations/19990054671](https://ntrs.nasa.gov/citations/19990054671)
- 814 [38] Giuseppe Penna, Daniele Magazzeni, Fabio Mercorio, and Intrigila  
815 Benedetto. 2009. UPMurphi: A tool for universal planning on PDDL+  
816 problems. In *Proc. of the 19th Int. Conf. on Automated Planning and*  
817 *Scheduling (ICAPS’09). September 19-23, Thessaloniki, Greece*.
- 818 [39] AIPlan4EU H2020 European Project. 2021. [https://www.](https://www.aiplan4eu-project.eu)  
819 [aiplan4eu-project.eu](https://www.aiplan4eu-project.eu)
- 820 [40] Franco Raimondi, Charles Pecheur, and Guillaume Brat. 2007. Testing  
821 Planning Domains (without Model Checkers). *Electron. Notes Theor. Com-*  
822 *put. Sci.* 190, 2 (2007).
- 823 [41] Franco Raimondi, Charles Pecheur, and Guillaume Brat. 2009. PDVer,  
824 a Tool to Verify PDDL Planning Domains. In *Proc. of ICAPS’09 Work-*

- 825 *shop on Verification and Validation of Planning and Scheduling Systems.*  
826 *September 19-23, Thessaloniki, Greece.*
- 827 [42] Robert Sasak. 2021. Prolog Planning Library. [https://github.com/  
828 RobertSasak/Prolog-Planning-Library](https://github.com/RobertSasak/Prolog-Planning-Library).
- 829 [43] Sergio Segura, Gordon Fraser, Ana B. Sanchez, and Antonio Ruiz-Cortes.  
830 2016. A Survey on Metamorphic Testing. *IEEE Transactions on Soft-  
831 ware Engineering* 42, 9 (2016), 805–824. [https://doi.org/10.1109/  
832 TSE.2016.2532875](https://doi.org/10.1109/TSE.2016.2532875)
- 833 [44] Sergio Segura, Dave Towey, Zhi Quan Zhou, and Tsong Yueh Chen. 2020.  
834 Metamorphic Testing: Testing the Untestable. *IEEE Software* 37, 3 (2020),  
835 46–53. <https://doi.org/10.1109/MS.2018.2875968>
- 836 [45] Anas Shrinah and Kerstin Eder. 2020. Goal-constrained planning domain  
837 model verification of safety properties. In *Proceedings of the 9th European  
838 Starting AI Researchers' Symposium 2020 co-located with 24th European  
839 Conference on Artificial Intelligence (ECAI 2020) (CEUR Workshop Pro-  
840 ceedings, Vol. 2655)*.
- 841 [46] Margaret H Smith, Gerard J Holzmann, Gordon C Cucullu, and BD Smith.  
842 2005. Model Checking Autonomous Planners: Even the Best Laid Plans  
843 Must be Verified. In *Proc. of the IEEE Aerospace Conference. Big Sky,  
844 Montana. In March*.
- 845 [47] Helge Spieker and Arnaud Gotlieb. 2020. Adaptive metamorphic testing  
846 with contextual bandits. *Journal of Systems and Software* 165 (2020).  
847 <https://doi.org/10.1016/j.jss.2020.110574>
- 848 [48] Thierry Titchou Chekam, Mike Papadakis, Tegawendé F. Bissyandé, Yves  
849 Le Traon, and Koushik Sen. 2020. Selecting Fault Revealing Mutants.  
850 *Empirical Software Engineering* 25, 1 (2020), 434–487. [https://doi.  
851 org/10.1007/s10664-019-09778-7](https://doi.org/10.1007/s10664-019-09778-7)

- 852 [49] Dongwei Xiao, Zhibo LIU, Yuanyuan Yuan, Qi Pang, and Shuai Wang.  
853 2022. Metamorphic Testing of Deep Learning Compilers. In *Proceedings*  
854 *of the ACM on Measurement and Analysis of Computing Systems*, Vol. 6.  
855 Association for Computing Machinery, New York, NY, USA, Article 15,  
856 28 pages. <https://doi.org/10.1145/3508035>
- 857 [50] Xiaoyuan Xie, Joshua Ho, Christian Murphy, Gail Kaiser, Baowen Xu,  
858 and Tson Chen. 2009. Application of Metamorphic Testing to Supervised  
859 Classifiers. In *2009 Ninth Int. Conf. on Quality Software*, Vol. 33. 135–144.
- 860 [51] Xiaoyuan Xie, Joshua W.K. Ho, Christian Murphy, Gail Kaiser, Baowen  
861 Xu, and Tsong Yueh Chen. 2011. Testing and Validating Machine Learning  
862 Classifiers by Metamorphic Testing. *Journal of Systems and Software* 84,  
863 4 (2011), 544–558. <https://doi.org/10.1016/j.jss.2010.11.920>
- 864 [52] Liming Xu, Dave Towey, Andrew P. French, Steve Benford, Zhi Quan Zhou,  
865 and Tsong Yueh Chen. 2018. Enhancing Supervised Classifications with  
866 Metamorphic Relations. *Proceedings of the 3rd International Workshop on*  
867 *Metamorphic Testing - MET '18* (2018), 46–53. [https://doi.org/10.](https://doi.org/10.1145/3193977.3193978)  
868 [1145/3193977.3193978](https://doi.org/10.1145/3193977.3193978)
- 869 [53] Shin Yoo. 2010. Metamorphic Testing of Stochastic Optimisation. In *2010*  
870 *Third International Conference on Software Testing, Verification, and Val-*  
871 *idation Workshops*. 192–201. <https://doi.org/10.1109/ICSTW.2010.26>
- 872 [54] Jiantao Zhang, Zheng Zheng, Beibei Yin, Kun Qiu, and Yang Liu. 2019.  
873 Testing Graph Searching Based Path Planning Algorithms by Metamor-  
874 phic Testing. In *2019 IEEE 24th Pacific Rim International Symposium on*  
875 *Dependable Computing (PRDC)*. 158–15809. [https://doi.org/10.1109/](https://doi.org/10.1109/PRDC47002.2019.00046)  
876 [PRDC47002.2019.00046](https://doi.org/10.1109/PRDC47002.2019.00046)
- 877 [55] Xiao-Yi Zhang, Yang Liu, Paolo Arcaini, Mingyue Jiang, and Zheng Zheng.  
878 2024. MET-MAPF: A Metamorphic Testing Approach for Multi-Agent  
879 Path Finding Algorithms. *ACM Trans. Softw. Eng. Methodol.* (jun 2024).  
880 <https://doi.org/10.1145/3669663> Just Accepted.

881 [56] Zhi Quan Zhou and Liqun Sun. 2019. Metamorphic Testing of Driverless  
882 Cars. *Commun. ACM* 62, 3 (2019), 61–67. <https://doi.org/10.1145/>  
883 [3241979](https://doi.org/10.1145/3241979)

Definition	Description
$f_1(s) = h(s)$	Expands the best states first.
$f_2(s) = -h(s)$	Expands the worst states first.
$f_3(s) = -g(s) + h(s)$	Encourages the search to expand the deepest nodes. It uses $h(s)$ to mitigate its naive trend to always search for a longer solution, thus leading to a timeout signal.
$f_4(s) = (-1)^{g(s)}h(s)$	Scores the discovered state with $h$ but “discards” it by signing it with the parity of its cost.
$f_5(s) = -g(s) + (-1)^{g(s)}h(s)$	Combination of $f_3$ and $f_4$ .
$h_1(s) =  s $	Ranks a state by its length (i.e, its number of true facts).
$h_2(s) =  g \setminus (s \cap g) $	Ranks a state with the number of facts of the goal that are not true yet.
$h_3(s) = s\Delta i$	Ranks a state with its distance to the initial state. The distance used is the symmetric difference.
$h_4(s) = s\Delta g$	Ranks a state with its distance to the goal state.
$h_5(s) =  succ(s) $	Ranks a state with its number of applicable actions. If we suppose that every action leads to a distinct state, then $h_5$ can be defined as the number of successors of the current state.
$h_6(s) = h_{max}$	Ranks a state with the index of the first fact layer of the relaxed planning graph, which contains all the facts of the goal. In our planning setting, where all action costs equal 1, it corresponds to the classical heuristic $h_{max}$ .

Table 2: Description of the  $f$  and  $h$  mutated functions used in a flawed  $A^*$  to create the mutated AI planners.  $g(s)$  refers to the cost of the plan used to reach the state  $s$ .  $g$  and  $i$  refer to the goal and the initial state, respectively.

	$f_1(s)$	$f_2(s)$	$f_3(s)$	$f_4(s)$	$f_5(s)$
$h_1(s)$	✓	✓	✓	✓	✓
$h_2(s)$	✓	✓	✓	E	E
$h_3(s)$	E	✓	✓	✓	✓
$h_4(s)$	E	✓	✓	E	✓
$h_5(s)$	✓	✓	✓	✓	✓
$h_6(s)$	E	✓	✓	E	E

Table 3: Overview of mutant combinations of  $f$  and  $h$ : ✓ indicates the used mutants, E marks equivalent mutants (unused).

	MORPHINPLAN		bfs_det		bfs_ran		wal_ran	
	Mut. Score	t (s)	Mut. Score	t (s)	Mut. Score	t (s)	Mut. Score	t (s)
airport06	<b>83.5±7.7</b>	175.8	25.3±7.5	17.9	79.6±5.4**	18.0	60.8±18.7***	2.8
airport07	<b>83.1±6.8</b>	176.2	13.6±5.9	13.3	75.3±8.2**	13.3	59.9±11.4***	1.6
blocks01	<b>82.2±6.9</b>	23.0	68.6±7.0	0.1	30.2±8.3*	0.1	18.2±7.9**	0.9
blocks02	<b>94.5±3.9</b>	23.0	66.5±8.7	0.1	74.6±7.2*	0.1	34.8±9.1**	1.2
blocks03	<b>83.9±8.9</b>	22.7	55.9±7.4	0.1	57.1±11.2**	0.1	32.8±8.9**	1.0
gripper01	<b>76.6±10.9</b>	29.9	33.0±9.6	0.1	70.3±5.8**	0.2	22.7±7.5**	2.4
elevator02	<b>92.5±8.2</b>	20.7	66.4±6.3	0.2	78.6±8.7*	0.1	3.9±4.6*	1.1
elevator03	<b>94.2±8.4</b>	29.4	55.7±8.6	0.1	70.9±7.1**	0.5	23.1±6.2**	2.1
openstacks01	<b>27.4±9.9</b>	213.9	12.7±8.1	0.8	17.2±6.1**	0.8	22.9±6.0**	0.2
pegsol04	<b>90.8±7.4</b>	45.8	72.6±7.8	3.4	45.2±15.5**	3.3	42.2±13.2**	1.0
pegsol05	<b>86.2±6.9</b>	51.3	33.6±7.7	4.7	46.5±24.1**	4.7	23.7±12.6***	1.1
pegsol06	<b>98.6±6.1</b>	147.5	29.1±8.5	6.9	45.0±12.7**	6.9	48.7±10.2***	1.6
pst-small03	<b>83.3±8.6</b>	22.7	41.8±9.1	0.1	41.3±11.0**	0.1	15.3±8.2**	0.5
pst-small06	<b>94.4±4.7</b>	27.4	53.5±5.7	0.2	76.0±9.7*	0.2	51.6±8.7**	0.9
pst-small08	<b>76.3±6.7</b>	27.5	56.1±7.9	0.2	48.8±10.1**	0.1	30.7±8.4**	0.8
pst-small09	<b>100.0±0.0</b>	26.9	84.0±5.6	0.2	57.5±10.2**	0.2	37.4±9.7**	1.1
tpp03	<b>75.8±9.4</b>	23.0	16.1±7.7	0.2	46.3±10.8**	0.2	18.2±8.5**	104.5
transport01	<b>96.0±5.3</b>	28.8	83.2±9.0	0.2	81.3±8.2*	0.3	49.9±9.1**	1.7
travel02	<b>89.5±4.4</b>	22.2	58.2±10.2	0.5	49.4±9.5**	0.5	49.6±5.9**	4.6
Mean±std. dev.	<b>84.7±15.3</b>	59.9	48.7±22.2	2.6	57.4±18.0	2.6	34.0±15.8	6.9

Table 4: Average mutation scores of MORPHINPLAN on the 19 problems for the 20 subsets of validation/selection mutants. For the non-deterministic methods, the standard deviation over 10 runs is indicated (\*  $\leq 5$ , \*\*  $\leq 10$ , \*\*\*  $\leq 15$ ).

Sasak's PPL [42]		FastDownward														
V0		V1		Heuristic (Non-Optimal) Configurations										Opt.		
	$A_{h0}$	$A_{hmax}$	$A_{ff}$	$A_{add,cca,cg}$	$W_{ff}$	$W_{add,cca}$	$W_{cg}$	$W_{goalcount}$								
blocks01	✗	✗	✗	✗	✗	✓	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗
blocks02	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
blocks03	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
elevator02	✓	✗	✗	✓	✗	✓	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗
elevator03	✓	✗	✗	✓	✓	✓	✓	✗	✗	✗	✗	✓	✓	✗	✗	✗
pstr-small03	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
pstr-small06	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
pstr-small08	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
pstr-small09	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
airport06	N/A	✗	N/A	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
airport07	N/A	✗	N/A	✗	✗	✓	✓	✗	✗	✗	✓	✗	✗	✗	✗	✗
pegso04	N/A	✗	N/A	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
pegso05	N/A	✗	N/A	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
pegso06	N/A	✗	N/A	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
gripper01	N/A	N/A	N/A	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
openstacks01	N/A	N/A	N/A	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
tpp03	N/A	N/A	N/A	✗	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓
transport01	N/A	N/A	N/A	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓
travel02	N/A	N/A	N/A	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗

Table 5: MORPHINPLAN results of non-optimal plans detection on 20 configs of two off-the-shelf AI planners. Columns showing similar behaviour have been merged. Checks indicate problems which have led to at least one non-optimal plan. Crosses indicate that no such plans were detected. N/A are cases where the planner failed due to unsupported PDDL features.