

# Formal Explanations of Neural Network Policies for Planning

This paper is under review for IJCAI. Due to the space limit of the HAXP submission, we are unable to take into account the comments of the IJCAI reviewers.

## Abstract

1 Deep learning is increasingly used to learn policies for plan- 42  
2 ning problems. However, policies represented by neural net- 43  
3 works are difficult to interpret, verify and trust. Existing formal 44  
4 approaches to post-hoc explanations provide concise reasons 45  
5 for a *single* decision made by an ML model. However, under- 46  
6 standing planning policies requires explaining *sequences* of 47  
7 decisions. In this paper, we formulate the problem of finding 48  
8 explanations for the sequence of decisions recommended by 49  
9 a learnt policy in a given state. We show that, under certain 50  
10 assumptions, a minimal explanation for a sequence can be 51  
11 computed by solving a number of single decision explanation 52  
12 problems which is linear in the length of the sequence. We 53  
13 present experimental results of our implementation of this 54  
14 approach for ASNNets policies for classical planning domains. 55

## 1 Motivation

15  
16 Deep learning has become the method of choice in the areas 60  
17 of AI that focus on perception, and is rapidly gaining traction 61  
18 in other areas that have traditionally been strongholds of rea- 62  
19 soning, search, and combinatorial optimisation. In automated 63  
20 planning for instance, new work has emerged that uses deep 64  
21 learning to learn policies and heuristics in a wide range of 65  
22 planning domains. We refer the reader to (Toyer et al. 2018; 66  
23 Groshev et al. 2018; Garg, Bajpai, and Mausam 2020; Zhang 67  
24 and Geißer 2022; Karia and Srivastava 2022) for examples 68  
25 of work aiming at learning policies for planning domains, 69  
26 and to (Shen, Trevizan, and Thiébaux 2020; Ferber, Helmert, 70  
27 and Hoffmann 2020; Karia and Srivastava 2021; Ferber et al. 71  
28 2022; Gehring et al. 2022) as representatives of work on 72  
29 learning heuristics to guide the search for a plan. 73

30 As the use of deep learning becomes more widespread in 74  
31 planning, the need to understand the solutions it produces 75  
32 becomes more pressing. Policies represented by neural net- 76  
33 works are notoriously opaque, and difficult to understand, 77  
34 verify, and trust (Toyer et al. 2020; Vinzent, Steinmetz, and 78  
35 Hoffmann 2022). At the minimum, one would like to be 79  
36 able to explain why a particular course of action was recom- 80  
37 mended by the policy – by identifying the properties of the 81  
38 state of the world prior to the execution of the policy which 82  
39 led to that recommendation – so as to help the user decide 83  
40 whether this recommendation should be trusted. This is the 84  
41 problem addressed in this paper. 85

## 1.1 Running Example

Throughout the paper, we use the example of Yvette who 42  
needs to take a turnpike to get to her final destination where 43  
she will spend a couple of weeks holidays. The turnpike 44  
requires to either purchase a weekly pass online or pay with 45  
cash at a toll gate. The pass is expensive and should not be 46  
taken unless one expects to use the turnpike multiple times in 47  
a single week. In our example, we assume that Yvette does 48  
not currently have a pass or cash. Hence, the policy prescribes 49  
to drive to an ATM, withdraw some money, drive to the toll 50  
gate, pay the toll, and drive to the destination. 51

52 Why choose this course of action? Firstly because i) Yvette 53  
is on the side of the turnpike opposite to her destination. 54  
Indeed this course of action would work in any state in which 55  
condition i) holds. However, this fails to explain *why* Yvette 56  
should not directly go to the toll gate. This is because ii) she 57  
has no pass and iii) she has no money. A correct explanation 58  
should therefore include all three conditions. 59

60 If, instead, Yvette did hold a pass, then the policy would 61  
skip the visit to the ATM. The explanation would then men- 62  
tion the fact that Yvette holds a pass but would not mention 63  
her lack of money as the policy would still have prescribed 64  
this course of action even if she had money. 65

66 Explanations can also expose unexpected reasons for deci- 67  
sions. Policies learnt using techniques such as deep learning 68  
are not guaranteed to be rational or even valid. They can 69  
also expose preferences or bias. For instance, the explanation 70  
could include the proposition iv) it is sunny, which implicitly 71  
means that the policy would have decided differently if it was 72  
not. It is questionable whether iv) is relevant in this context 73  
(maybe the idea is that you would not want to make a trip to 74  
the ATM under the rain). The neural network could also be 75  
prejudiced against certain groups of people and give different 76  
advice depending on gender, race, etc. (Darwiche and Hirth 77  
2020). 78

## 1.2 Existing Work

79 The nascent work on explainable planning has been focusing 80  
on a different set of problems in a different setting, in partic- 81  
ular on model reconciliation and contrastive explanations for 82  
conventional model-based planning (Chakraborti, Sreedharan, 83  
and Kambhampati 2020). In model reconciliation, the 84  
aim is to generate explanations allowing a human user to 85  
update his model of the planning problem to make it consis- 86  
tent with the plan produced by a planning agent (Chakraborti 87  
et al. 2017; Sreedharan, Chakraborti, and Kambhampati 2021; 88  
Vasileiou et al. 2022). This latter question is concerned with 89

the properties of the planning model, rather than those of the policy. The second prominent line of research in the explainable planning literature is the generation of contrastive explanations outlining why the planner chose a course of action over others within the space of possible plans (Eifler et al. 2020; Kasenberg, Thielstrom, and Scheutz 2020; Krarup et al. 2021). These works are concerned with understanding the space of possible decisions and their respective merits, rather than a particular policy.

Therefore, as a starting point, we instead turn to prior work concerned with explaining deep learning and other data-driven models for *classification* tasks. Existing approaches typically either compute simpler models that locally approximate the classifier’s behavior (Ribeiro, Singh, and Guestrin 2016; Lundberg and Lee 2017), or identify sufficient conditions on the inputs that led the neural network to produce a particular output (Ribeiro, Singh, and Guestrin 2018; Ignatiev, Narodytska, and Marques-Silva 2019). One approach falling into the latter class is to compute *abductive* explanations that are *minimal* sufficient conditions for the decision. This has the advantage of providing formal guarantees of soundness and non-redundancy (Ignatiev, Narodytska, and Marques-Silva 2019; Darwiche and Hirth 2020; Marques-Silva and Ignatiev 2022).

However, the above approaches are designed to explain a single decision, whereas understanding the recommendations of a planning policy requires explaining why a particular *sequence* of decisions was made. The latter is more challenging as it involves reasoning about repeated applications of the policy network and about the successive changes they induce in the state of the world in which the policy is executed.

### 1.3 Contribution

In this paper, we extend abductive explanations from single to sequential decisions. We restrict ourselves to classical planning policies and explanations of why the policy makes a certain sequence of decisions from a given state. We formally define explanations for a sequence of decisions, and show that, under certain assumptions, the problem of finding an explanation for the sequence can be decomposed into that of finding explanations for the individual decisions in the sequence. We provide an algorithm that exploits this decomposition to compute a minimal explanation for a sequence by making a number of consistency tests pertaining to individual decisions that is at most linear in the length of the sequence and in the number of state variables. We then discuss the implementation of our approach to explain policies represented by Action Schema Networks (ASNs) (Toyer et al. 2020) and report on its performance on sparse ASNs policies for classical planning domains. We conclude by discussing the limits of our work and possible extensions.

## 2 Background

We start by introducing the type of planning problems we consider, their representation, and our notations.

Many of the recent work on deep learning for planning assume that the model of the planning domain is available to the learner. We assume that it is also available to the explainer. Here we represent the classical planning instance

$I = \langle X, A, g \rangle$  under consideration using the SAS<sup>+</sup> formalism (Bäckström and Nebel 1995).  $X$  is a set of finite-domain *state variables*, where  $D_x$  is the *domain* of variable  $x$ . A *partial state* (or partial valuation)  $s$  is an assignment of value to a subset  $X_s \subseteq X$  of the variables such that  $s[x] \in D_x$  for  $x \in X_s$ . If  $X_s = X$  then we say that  $s$  is a *state*; we write  $S$  for the set of states. A *proposition* ( $x = v$ ) is a partial valuation assigning a value to a single variable.

The goal  $g$  is a partial state. Given two partial states  $s$  and  $s'$ , we write  $s \subseteq s'$  when  $s[x] = s'[x]$  for all  $x \in X_s$ . A *completion* of a partial state  $s$  is a state  $s'$  such that  $s \subseteq s'$ . The result of applying a partial valuation  $e$  to a partial state  $s$  is the partial state  $s \oplus e$  over  $X_s \cup X_e$  defined by  $(s \oplus e)[x] = e[x]$  if  $x \in X_e$  and  $(s \oplus e)[x] = s[x]$  if  $x \in X_s \setminus X_e$ . We also define the binary operator  $\ominus$  over partial states:  $s \ominus e$  is the restriction of  $s$  to the variables  $X_s \setminus X_e$ . For a variable  $x \in X_s$  will write  $s - x$  as an abbreviation for  $s \ominus (x = s[x])$ .

$A$  is the set of *actions*. Action  $a \in A$  is characterised by two partial valuations representing its *precondition*  $\text{pre}(a)$  and its *effect*  $\text{eff}(a)$ , respectively. We say that the action is *applicable* in a state  $s \in S$  iff  $\text{pre}(a) \subseteq s$  and write  $A(s) \subseteq A$  for the subset of actions applicable in  $s$ . Moreover, given a partial state  $s$  and an action  $a$ , the *progression* of  $s$  through  $a$  is the partial state  $\text{prg}_a(s) = s \oplus \text{eff}(a)$ . Note that if  $s$  is a state and  $a \in A(s)$  then  $\text{prg}_a(s)$  is the state resulting from applying  $a$  in  $s$ .

A *policy* for the planning instance is a function  $\pi : S \mapsto A$  mapping states to applicable actions, i.e.  $\pi(s) \in A(s)$ . We define the *n-long trajectory*  $\tau_\pi^n(s)$  induced by  $\pi$  from state  $s$  as follows:  $\tau_\pi^n(s) = s_1 \xrightarrow{a_1} \dots \xrightarrow{a_n} s_{n+1}$  such that  $s_1 = s$  and for all  $0 < i \leq n$ ,  $a_i = \pi(s_i)$  and  $s_{i+1} = \text{prg}_{a_i}(s_i)$ . Finally, the *n-long sequence of actions recommended* by  $\pi$  in  $s$  is  $\pi^n(s) = a_1, \dots, a_n$ , where the  $a_i$ s are the successive actions in  $\tau_\pi^n(s)$ .

Given a planning instance  $I$ , a policy  $\pi$  for  $I$ , an *initial state*  $s$ , and an integer  $n$ , our problem is to explain why  $\pi$  recommended the sequence of actions  $\pi^n(s)$  in  $s$ . The explanation is meant to shed light on the appropriateness of the recommendation.

The definition of explanations presented in the next section relies on the fact that the policy will recommend the same sequence of actions for certain states. To ensure that for any state  $s$  and any length  $n$ , a policy can recommend an  $n$ -long sequence from  $s$ , we make the following assumptions. First, we assume that there is no terminal state, i.e.,  $A(s) \neq \emptyset$  for all states. This is not a restriction as a default action could be to do nothing. Second, we assume that  $A$  includes a special goal action  $a_g$  which has no effect and which is only applicable in goal states (this will act as a marker of the goal being reached), and that  $\pi$  is a total function on  $S$  such that  $\pi(s) = a_g$  iff  $g \subseteq s$ . This is purely for convenience as policies generally check whether a goal state has been reached before computing the next action. With these assumptions, our theory applies uniformly to any trajectory, regardless of whether it reaches the goal.

## 3 Explanations of Neural-Network Policies

An explanation of a decision (or sequence of decisions) in a state is a condition on this state that led to this decision

203 being made: the decision was made because the condition was  
 204 satisfied in this state. Said differently, the same decision would  
 205 have been made in any other state that satisfies this condition.  
 206 We could allow arbitrary conditions; e.g. the explanation  
 207 could be the logical formula that describes exactly all the  
 208 states in which this decision would be taken. However, such  
 209 an explanation would not be very helpful. We aim instead for  
 210 a ‘simple’ explanation, that is, an explanation that mentions  
 211 as few propositions as possible and has the simple structure  
 212 of a conjunction.

213 We use a definition similar to that of (Marques-Silva and  
 214 Ignatiev 2022). An explanation is a partial state that entails  
 215 the decision; in logic, this is known as an implicant.

216 **Definition 1.** An explanation of a single decision  $a$  for policy  
 217  $\pi$  is a partial state  $\mathbf{z}$  such that  $\pi$  yields the same decision for  
 218 all completions of  $\mathbf{z}$ :

$$\forall s \in S. (\mathbf{z} \subseteq s) \implies \pi(s) = a.$$

219 When  $s$  completes  $\mathbf{z}$ , we say that  $\mathbf{z}$  explains decision  $a$  in  $s$ .

220 Our goal is not to explain just the first decision of the  
 221 policy, but the complete sequence of decisions. While the first  
 222 decision was based on the initial state, later decisions were  
 223 made based on the later states. These states, however, are fully  
 224 determined by the initial state and the actions taken. Using the  
 225 planning model, it is therefore possible to trace the sufficient  
 226 condition that led to the full sequence of actions back to the  
 227 initial state.

228 **Definition 2.** An explanation of the  $n$ -long sequence of deci-  
 229 sions  $a_1, \dots, a_n$  for a policy  $\pi$  is a partial state  $\mathbf{z}$  such that  
 230  $\pi$  yields the same sequence of decisions for all completions  
 231 of  $\mathbf{z}$ :

$$\forall s \in S. (\mathbf{z} \subseteq s) \implies \pi^n(s) = a_1, \dots, a_n.$$

232 Similarly as before, when  $s$  completes  $\mathbf{z}$ , we say that  $\mathbf{z}$  ex-  
 233 plains the sequence of decisions  $a_1, \dots, a_n$  in  $s$ . We note that  
 234 if  $a_n = a_g$  is the goal action, then the sequence  $a_1, \dots, a_{n-1}$   
 235 leads all completions of  $\mathbf{z}$  to a goal state since  $\pi$  only recom-  
 236 mends applicable actions and  $a_g$  is only applicable in a goal  
 237 state. The sequence of actions recommended by the policy  
 238 might not lead to the goal; in this case, the loop of the infi-  
 239 nite trajectory induced by the policy does not occur at a goal  
 240 state. If one wants to compute an explanation for this infinite  
 241 sequence, it is possible to use Definition 2 with  $n = L \times |S|$   
 242 where  $L$  is the length of the loop and  $|S|$  the total number of  
 243 states: if  $\mathbf{z}$  explains  $\pi^n(s)$ , then it explains  $\pi^{n+k}(s)$  for all  
 244  $k \geq 0$ . It may be possible to derive better bounds.

245 It should be clear that there can be multiple explanations in  
 246 the same state. For instance, in our running example, Yvette  
 247 needs either a pass or some cash to take the turnpike. Both the  
 248 fact that she has a pass and the fact that she has cash would be  
 249 acceptable explanations for driving directly to the toll gate. In  
 250 a state where she has a pass and cash, these two explanations  
 251 are therefore suitable. We also note that explanations enjoy  
 252 monotonic properties: if  $\mathbf{z}$  is an explanation, any superset of  
 253  $\mathbf{z}$  is an explanation. In particular, the complete initial state is  
 254 an explanation, although hardly a useful one.

255 Our goal is to compute the ‘best’ explanation for the se-  
 256 quence of decisions made from our initial state. Specifically,

257 we want to compute a subset-minimal explanation (akin to  
 258 a *prime* implicant in logic), i.e., an explanation  $\mathbf{z}$  such that  
 259 no strict subset of  $\mathbf{z}$  is an explanation. Minimal explanations  
 260 provide additional benefits: all variables mentioned in the ex-  
 261 planation are required, in the sense that if any were removed,  
 262 the partial state would no longer be an explanation. Therefore,  
 263 seeing a variable that should not be relevant in a minimal ex-  
 264 planation should raise questions about the policy, while this  
 265 phenomenon is unsurprising in a non-minimal explanation.

266 **Definition 3.** Given a policy  $\pi$ , an integer  $n$ , and a state  $s$ ,  
 267 the minimal explanation problem is to find a minimal partial  
 268 state that explains the sequence of decisions  $\pi^n(s)$  in  $s$ .

269 Ignatiev, Narodytska, and Marques-Silva (2019) have  
 270 shown how to compute explanations for single decisions. The  
 271 policy is translated into a set of constraints  $C_\pi$  over a set of  
 272 variables that includes the state variables  $X$  which are the  
 273 input to the policy, and the variable  $y$  which represents its  
 274 output. The model of  $C_\pi$  are exactly all the pairs  $\langle s, y \rangle, s \in S$ ,  
 275  $y = \pi(s)$ . If  $\pi$  is represented by a neural network,  $C_\pi$  can  
 276 be formulated as a set of mixed-integer programming or sat  
 277 modulo theory constraints (see Section 5). In order to decide  
 278 whether  $\mathbf{z}$  is an explanation for a decision  $a$ , the constraints  $\mathbf{z}$   
 279 and  $\neg d$  are added to the set where  $d \equiv (y = a)$ . If the result-  
 280 ing set of constraints is consistent, then there exists a state  $s'$   
 281 that completes  $\mathbf{z}$  and yields a decision different from  $a$ ; hence,  
 282  $\mathbf{z}$  does not explain  $a$ . Otherwise, all states that complete  $\mathbf{z}$  lead  
 283 to decision  $a$ , and  $\mathbf{z}$  explains  $a$ :

$$\mathbf{z} \text{ explains } a \Leftrightarrow \neg \text{CO}(C_\pi \wedge \mathbf{z} \wedge \neg d).$$

284 Using monotonicity, it is then possible to greedily search  
 285 for a minimal explanation. This is done by starting with an  
 286 existing explanation  $\mathbf{z}$ , for instance the initial state  $s$ , and  
 287 testing whether for some variable  $x \in X_{\mathbf{z}}$ ,  $\mathbf{z}' := \mathbf{z} - x$   
 288 remains an explanation. If  $\mathbf{z}'$  indeed explains  $a$ , we replace  
 289  $\mathbf{z}$  with it; otherwise we move to the next variable  $x$  until we  
 290 tried to remove each variable. The same variable does not  
 291 need to be tested more than once.

## 292 4 Computing a Minimal Explanation

293 We now turn to the problem of computing a minimal expla-  
 294 nation for a sequence of decisions.

### 295 4.1 Naive Algorithm

296 For completeness, we first consider a naive algorithm, il-  
 297 lustrated in Figure 1. We expect that this algorithm will be  
 298 impractical, as it requires testing the consistency of constraint  
 299 sets involving too many variables.

300 Similarly as in the single decision case, the idea of the  
 301 algorithm is to build a single set of constraints which is con-  
 302 sistent iff a specified partial state does not explain a specified  
 303 sequence of decisions  $a_1, \dots, a_n$ . Given an initial state  $s_1$ ,  
 304 we define the set of constraints  $C_{a_1, \dots, a_n}$  which computes the  
 305 states  $s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots \xrightarrow{a_{n-1}} s_n$  reached by applying the  
 306 successive actions as well as the decisions  $y_i = \pi(s_i)$  of the  
 307 policy in each of these states; we then compare the  $y_i$ s with

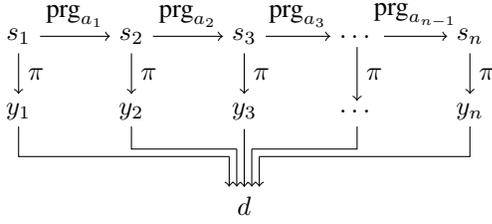


Figure 1: Graphical representation of  $C_{a_1, \dots, a_n}$ , the set of constraints used to determine whether a partial state is an explanation. Nodes of the graph are sets of variables. Arrows represent constraints defined such that the target variables are a function of the source variables.

308 the  $a_i$ :

$$\begin{aligned} s_{i+1} &= \text{prg}_{a_i}(s_i) & \forall i \in \{1, \dots, n-1\} \\ y_i &= \pi(s_i) & \forall i \in \{1, \dots, n\} \\ d_i &= (y_i = a_i) & \forall i \in \{1, \dots, n\} \\ d &= \bigwedge_{i=1}^n d_i \end{aligned}$$

309 Finally, we add the constraints that  $s_1$  should complete  $\mathbf{z}$  and  
310 that  $d$  should be false. The resulting set of constraints,

$$C_{a_1, \dots, a_n} \wedge \bigwedge_{x \in X_{\mathbf{z}}} (s_1[x] = \mathbf{z}[x]) \wedge \neg d,$$

311 is consistent iff there exists a state  $s_1$  that completes  $\mathbf{z}$  for  
312 which  $\pi$  generates a different sequence of decisions than  
313 the specified one; i.e.,  $\mathbf{z}$  does not explain  $a_1, \dots, a_n$ . As  
314 before, one can then use a greedy algorithm to compute the  
315 explanation, starting with  $\mathbf{z} = s_1$  and progressively removing  
316 propositions.

317 However in this set of constraints, the variables include  $n$   
318 duplicates of the set of state variables  $X$ , which are needed to  
319 represent the successive states  $s_1 \dots s_n$ , as well as the much  
320 larger set of variables required to represent  $n$  duplicates of the  
321 computation of  $\pi$  (via the constraints  $C_\pi$  mentioned above).  
322 This makes testing consistency impractical for anything but  
323 very small sequences.

## 324 4.2 Forward Decomposition

325 We now show that it is possible to decompose the explanation  
326 so that it is easier to compute a single minimal explanation.  
327 This decomposition leads to an algorithm that need only solve  
328 a number of single decision explanation problems that is in  
329 the worst case linear in the length of the sequence. In the  
330 following  $\circ$  is the function composition.

331 **Theorem 1.** *Let  $\pi$  be a policy,  $s$  be a state, and let  $\tau_\pi^n(s) =$   
332  $s_1 \xrightarrow{a_1} \dots \xrightarrow{a_n} s_{n+1}$  be the  $n$ -long trajectory induced by  $\pi$   
333 from  $s$ . Let  $\mathbf{z} \subseteq s$  be a partial state and let  $\mathbf{z}_1, \dots, \mathbf{z}_n$  be a  
334 sequence of partial states defined by  $\mathbf{z}_{i+1} = \text{prg}_{a_i} \circ \dots \circ$   
335  $\text{prg}_{a_1}(\mathbf{z})$ . Then  $\mathbf{z}$  explains  $\pi^n(s)$  iff for all steps  $i$ ,  $\mathbf{z}_i$  explains  
336  $a_i$  in  $s_i$ .*

337 The proof of Theorem 1 is in Appendix A. Theorem 1 gives  
338 us a clear procedure to verify whether a partial state  $\mathbf{z}$  is an  
339 explanation, which is to compute the partial states  $\mathbf{z}_i$  resulting  
340 from applying the actions from  $\mathbf{z}$  (note:  $\mathbf{z}_{i+1} = \text{prg}_{a_i}(\mathbf{z}_i)$ )  
341 and to verify whether each  $\mathbf{z}_i$  explains  $a_i$ .

Algorithm 1: Computing a minimal explanation for a sequence of decisions.

---

```

1: procedure MINIMALEXPLANATION( $I, \pi, n, s$ )
2:    $s_1 \xrightarrow{a_1} \dots \xrightarrow{a_n} s_{n+1} = \tau_\pi^n(s)$ 
3:    $\mathbf{z}_1 := s$ 
4:   for  $i \in \{1, \dots, n-1\}$  do
5:      $\mathbf{z}_{i+1} := \text{prg}_{a_i}(\mathbf{z}_i)$ 
6:   for  $x \in X$  do
7:      $\text{explains} := \text{True}$ 
8:     for  $i \in \{1, \dots, n\}$  do  $\triangleright$  Testing condition of Th. 1
9:        $\mathbf{z}_i := \mathbf{z}_i - x$ 
10:      if  $\text{CO}(C_\pi \wedge \mathbf{z}_i \wedge \neg d_i)$  then
11:         $\text{explains} := \text{False}$ 
12:      break
13:      if  $x \in X_{\text{eff}(a_i)}$  then
14:        break  $\triangleright$  Can stop now (cf. Eq. 1)
15:      if  $\neg \text{explains}$  then  $\triangleright$  Must reinsert  $x$ 
16:        for  $j \in \{1, \dots, i\}$  do  $\mathbf{z}_j := \mathbf{z}_j \oplus (x = s_j[x])$ 
17:  return  $\mathbf{z}_1$ 

```

---

In the greedy algorithm that we propose next, we use the  
342 additional property: for any variable  $x$ ,  $\text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z} -$   
343  $x) =$   
344

$$\begin{cases} \text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z}) & \text{if } x \in X_{\text{eff}(a_i)} \cup \dots \cup X_{\text{eff}(a_1)} \\ \text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z}) - x & \text{otherwise.} \end{cases} \quad (1)$$

The first case is useful because it means  $\text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z} -$   
345  $x)$  explains  $a_{i+1}$  iff  $\text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z})$  does; so as soon  
346 as variable  $x$  appears in the effects of  $a_i$ , we will be able  
347 to ignore the condition of Theorem 1 that  $z_k$  should explain  
348  $a_k$  in  $s_k$  for all  $k > i$ . The second case is useful because it  
349 makes it easy to compute  $\text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z} - x)$  from  
350  $\text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z})$ .  
351

Our procedure is described in Algorithm 1. In Lines 3–5,  
352 the algorithm initialises the variables  $\mathbf{z}_i$ , mirroring the vari-  
353 ables from Theorem 1 for explanation  $\mathbf{z} := s$ . Line 6 starts  
354 the loop in which each variable  $x$  is tentatively removed from  
355 the explanation. The variable  $\text{explains}$  will record whether  
356  $\mathbf{z} - x$  explains the sequence; until proven otherwise, it is  
357 assumed it does. The inner loop from Line 8 verifies the con-  
358 dition of Theorem 1. After  $\mathbf{z}_i$  is updated (using the second  
359 case of Eq. 1), the condition is verified on Line 10. If  $\mathbf{z}_i$  does  
360 not explain action  $a_i$ , the inner loop is stopped; the loop from  
361 Line 16 will then undo the update. Otherwise, the inner loop  
362 moves to the next step  $i + 1$  *except* if variable  $x$  is mentioned  
363 in the effects of  $a_i$  in which case the loop can be stopped  
364 (first case of Eq. 1).  
365

This algorithm requires  $O(n \times m)$  calls to the consistency  
366 solver where  $n$  is the length of the sequence of  $m$  is the  
367 number of state variables.  
368

We illustrate this algorithm with our running exam-  
369 ple. We assume the policy described in Fig. 2 is used.  
370 The initial state, and so the initial explanation, is  $\mathbf{z} =$   
371  $\{\text{Friday}, \text{Sunny}, \text{NoMoney}, \text{NoPass}, \text{AtHome}, \dots\}$ .  
372

Algorithm 1 starts by removing *Friday* from the explana-  
373 tion. The partial state is therefore  $\mathbf{z}' = \mathbf{z}_1' = \mathbf{z} \ominus \{\text{Friday}\}$ .  
374

1. If has passed the TollGate, has a pass, or has some money, drive towards the destination.
2. Otherwise, if not at the ATM, drive towards the ATM.
3. Otherwise, if it is sunny, withdraw money.
4. Otherwise, purchase a pass online.

Figure 2: The policy given to Yvette (without goal action  $a_g$ )

We find that it is impossible to instantiate the only free variable (day of the week) in such a way that a decision different from *DriveToATM* is taken; in other words,  $\mathbf{z}'_1$  explains the first decision. Progressing  $\mathbf{z}'_1$  gives us  $\mathbf{z}'_2 = \{\text{Sunny}, \text{NoMoney}, \text{NoPass}, \text{AtATM}, \dots\}$ . This partial state also explains the second decision (*WithdrawMoney*), and so on until the end of the sequence. Therefore,  $\mathbf{z}'$  explains the whole sequence.

Then, Algorithm 1 removes *Sunny* from the explanation. The partial state is therefore  $\mathbf{z}'' = \mathbf{z}'_1 = \mathbf{z}' \ominus \{\text{Sunny}\}$ . Partial state  $\mathbf{z}''$  explains the first decision (*DriveToATM*). Progressing  $\mathbf{z}''_1$  gives us  $\mathbf{z}''_2 = \{\text{NoMoney}, \text{NoPass}, \text{AtATM}, \dots\}$ . This time, we find that a state in which Yvette is at the ATM and the weather is rainy will yield a different decision (*PurchasePass*) than the second decision (*WithdrawMoney*). Therefore,  $\mathbf{z}''$  does not explain the sequence and *Sunny* is added back to the explanation. Similarly, the algorithm would then try and fail to remove the 3 remaining propositions.

## 5 Implementation

We use Algorithm 1 to explain the recommendations of AS-Nets policies (Toyer et al. 2018) for classical planning domains. This requires implementing the consistency test in line 10 of the algorithm for ASNets, which have a complex structure and nonlinear activation functions. Our encoding, presented below, is supported by mixed integer programming (MIP) solvers such as Gurobi (Gurobi Optimization, LLC 2023).

### 5.1 ASNets

ASNets are neural networks dedicated to planning, trained to produce policies that apply to any problem instance from a given planning domain modelled in (P)PDDL (Younes and Littman 2004). An ASNet consists of  $L$  alternating *action* and *proposition layers*, beginning and ending with an action layer. In each action layer (resp. proposition layer), each action (resp. proposition) of the planning instance is represented by an action *module* (resp. proposition module). Modules in one layer are connected to *related* modules in the next layer, where a proposition  $p$  and action  $a$  are related, written  $R(a, p)$  iff  $p$  appears in the precondition or effect of  $a$ . This enables relevant information to be efficiently propagated through the network.

ASNets is capable of learning policies that generalise to problem instances of different size from the same domain thanks to a weight sharing mechanism whereby the action (resp. proposition) modules from the same layer that are ground instances of the same action schema (resp. the same predicate), have the same weights. Here we omit details such

as the use of skip connections, and the more complex definition of relatedness in (Toyer et al. 2020), but our implementation supports them.

**Action Layers** At layer  $l$ , excluding the 1st and last layer, the module for action  $a$  takes as input a vector  $u_a^l$  which is constructed by enumerating the propositions  $p_1, \dots, p_M$  that are related to  $a$  and concatenating the outputs  $\psi_{p_1}^{l-1}, \dots, \psi_{p_M}^{l-1}$  of these propositions' modules from the previous layer. That is  $u_a^l = [\psi_{p_1}^{l-1} \dots \psi_{p_M}^{l-1}]^T$ . The output of the module is the vector  $\phi_a^l = f(W_a^l u_a^l + b_a^l)$  where  $W_a^l$  is a weight matrix,  $b_a^l$  a bias vector, and  $f$  a nonlinearity. ASNets uses exponential linear units (ELU), i.e.  $f(x) = x$  if  $x \geq 0$  and  $e^x - 1$  otherwise. Note that except for the final layer, the output vectors of all modules in the network have the same dimension  $d$ .

**Proposition Layers** At layer  $l$ , the module for proposition  $p$  takes as input a vector  $v_p^l$  constructed by enumerating the actions that are related to the proposition, pooling from the outputs of their modules at layer  $l$  if they share the same action schema, and concatenating the results. That is  $v_p^l = [\text{pool}(\{\phi_a^l \mid \text{op}(a) = o_1 \wedge R(a, p)\}) \dots \text{pool}(\{\phi_a^l \mid \text{op}(a) = o_S \wedge R(a, p)\})]^T$ , where  $\text{pool}$  represents max-pooling and  $\text{op}(a)$  is the action schema of action  $a$ . Similarly as for action modules, the output of the proposition module is the vector  $\psi_p^l = f(W_p^l v_p^l + b_p^l)$ .

**First Layer** The input to an action module of the first layer are the boolean values (0 or 1) of all its related propositions in the current state, booleans indicating whether each of these propositions is true in the goal, and a boolean indicating whether the action is applicable. That is the input vector is  $u_a^1 = [\sigma \ \gamma \ \epsilon]^T$  where for all propositions  $p_1, \dots, p_M$  related to the action,  $\sigma \in \{0, 1\}^M$  with  $\sigma_j = 1$  iff  $p_j \in s$ ,  $\gamma \in \{0, 1\}^M$  with  $\gamma_j = 1$  iff  $p_j \in g$ , and  $\epsilon = 1$  iff  $a \in A(s)$ .

**Last Layer** The output of an action module in the final layer  $l = L$  is a single logit  $\phi_a^L$  representing the unnormalised probability that this action should be taken in the current state  $s$  given as input to the network. When, as in this paper, a deterministic policy is sought, the applicable action  $a \in A(s)$  with maximum  $\phi_a^L$  is selected by the policy.

### 5.2 MIP Encoding

The main purpose of the encoding is to answer consistency queries where only *some* of the inputs to the ASNets are given, and its output is constrained. The key decision variables in the MIP model fall into 3 categories: variables representing the network inputs, its outputs, and the action and proposition modules. As is well known from other MIP encodings of neural networks, one also needs auxiliary variables to encode the activation functions. Parameters include the weight matrices and bias vectors described above. The MIP has no objective function since we are only testing for consistency.

**Policy Inputs** We encode each element in the input of an ASNet (current state propositions, goal proposition, applicable actions) using the following binary variables:  $S_p$  is true iff proposition  $p$  is true in the current state,  $G_p$  is true iff proposition  $p$  is in the goal, and  $E_a$  is true iff  $a$  is applicable

477 in the current state. Since a key feature of the MIP model  
 478 is to allow specifying *partial* states as input, we must add  
 479 constraints capturing when actions are applicable: action  $a$  is  
 480 applicable when  $S_p = 1$  for all  $p \in \text{pre}(a)$ .

$$\begin{aligned} E_a &\leq S_p && \forall a \in A, \forall p \in \text{pre}(a) \\ E_a &\geq 1 - |\text{pre}(a)| + \sum_{p \in \text{pre}(a)} S_p && \forall a \in A \end{aligned}$$

481 Moreover, ASNet only supports *boolean* propositions  $p \equiv$   
 482  $(x = v)$  for  $x \in X$  and  $v \in D_x$ . To encode  $SAS^+$  state  
 483 variables, we add mutex constraints ensuring that at most one  
 484 proposition assigning a value to a given variable can be true.

$$\begin{aligned} \sum_{v \in D_x} S_{(x=v)} &\leq 1 && \forall x \in X \\ \sum_{v \in D_x} G_{(x=v)} &\leq 1 && \forall x \in X \end{aligned}$$

485 **Action Modules** To encode action modules, we define the  
 486 MIP variables  $(\text{PAC}_a^l)_i$  representing the  $i$ th element of the  
 487 pre-activation vector (omitting the non-linearity  $f$ ) of the  
 488 module for action  $a$  in layer  $l$ , and  $(\text{OUT}_a^l)_i$  representing  
 489 the  $i$ th element of the output vector  $\phi_a^l$  of this module. For  
 490 an action  $a$  with related propositions  $p_1, \dots, p_M$ , we have,  
 491 taking  $j = (m - 1)d + k$

$$\begin{aligned} (\text{PAC}_a^l)_i &= \sum_{m=1}^M \sum_{k=1}^d (W_a^l)_{i,j} \cdot (\text{OUT}_{p_m}^{l-1})_k + (b_a^l)_i \text{ for } l > 1 \\ (\text{PAC}_a^1)_i &= \sum_{m=1}^M (W_a^1)_{i,m} \cdot S_{p_m} + (W_a^1)_{i,M+m} \cdot G_{p_m} + \\ &\quad (W_a^1)_{i,2M+1} \cdot E_a + (b_a^1)_i \end{aligned}$$

492 **Proposition Modules** Similarly, we define  $(\text{PAC}_p^l)_i$  and  
 493  $(\text{OUT}_p^l)_i$  for each proposition  $p$ . For proposition modules we  
 494 additionally need variables  $(\text{POOL}_{p,o}^l)_i$  to represent the  $i$ th  
 495 element of the pooled vector over actions related to  $p$  sharing  
 496 the action schema  $o$ . If related actions belong to  $S$  action  
 497 schemas  $o_1, \dots, o_S$ , we have, taking  $j = (s - 1)d + k$

$$\begin{aligned} (\text{PAC}_p^l)_i &= \sum_{s=1}^S \sum_{k=1}^d (W_p^l)_{i,j} \cdot (\text{POOL}_{p,o_s}^l)_k + (b_p^l)_i \\ (\text{POOL}_{p,o}^l)_i &= \max(\{(\text{OUT}_a^l)_i \mid \text{op}(a) = o \wedge R(a, p)\}) \end{aligned}$$

498 **Activations** The encoding of  $f$  as the ELU function is very  
 499 similar to the encoding of ReLU in MIP (Fischetti and Jo  
 500 2018) and is the same for proposition and action modules.  
 501 Given a proposition or action  $b$  we define the auxiliary vari-  
 502 ables  $(t_b^l)_i$  and  $(s_b^l)_i$  to store the positive and negative com-  
 503 ponent of the variable  $(\text{PAC}_b^l)_i$ , respectively, and  $(z_b^l)_i$  which  
 504 is an indicator variable for the sign of  $(\text{PAC}_b^l)_i$ . This leads to  
 505 the following constraints (note that Gurobi linearises the ex-  
 506 ponential, see (Gurobi Optimization, LLC 2023, p584-586))

$$\begin{aligned} (z_b^l)_i = 1 &\implies (t_b^l)_i \leq 0, && (z_b^l)_i = 0 \implies (s_b^l)_i \leq 0 \\ &(t_b^l)_i \geq 0, && (s_b^l)_i \geq 0 \\ (\text{PAC}_b^l)_i &= (t_b^l)_i - (s_b^l)_i, && (\text{OUT}_b^l)_i = (t_b^l)_i + e^{-(s_b^l)_i} - 1 \end{aligned}$$

507 **Policy Output** The chosen action is the applicable action  
 508  $a$  maximising the output  $\phi_a^L$  of the final layer’s modules.  
 509 We enforce this using the following additional variables:  $C_a$   
 510 which is a binary variable true iff action  $a$  is chosen, and

$Vmax$  which is the maximal value of  $\phi_a^L$  for *applicable* 511  
 512 actions.

$$\begin{aligned} \sum_{a \in A} C_a &= 1 && \forall a \in A \\ C_a = 1 &\implies E_a = 1 && \forall a \in A \\ E_a = 1 &\implies Vmax \geq \text{OUT}_a^L && \forall a \in A \\ C_a = 1 &\implies Vmax \leq \text{OUT}_a^L && \forall a \in A \end{aligned}$$

**Temporary Constraints** The above constraints constitute 513  
 the encoding  $C_\pi$  of the ASNet policy and only needs to be 514  
 built once for a given sets  $X$  and  $A$  of state variables and 515  
 actions. For each consistency query involving a planning 516  
 instance  $I = \langle X, A, g \rangle$ , an action decision  $a$ , and a candidate 517  
 explanation  $\mathbf{z}$ , it suffices to temporarily add the following 518  
 constraints to set  $g$  as the goal, prevent  $a$  to be chosen by the 519  
 policy, and look for a possible completion  $s$  of  $\mathbf{z}$ . 520

$$\begin{aligned} C_a &= 0 \\ S_{(x=\mathbf{z}[x])} &= 1 && \forall x \in X_{\mathbf{z}} \\ G_{(x=g[x])} &= 1 && \forall x \in X_g \\ \sum_{v \in D_x} G_{(x=v)} &= 0 && \forall x \in X \setminus X_g \end{aligned}$$

## 6 Experimental Results 521

The goal of our experiments is to evaluate how effective 522  
 abductive explanations are in explaining why a policy recom- 523  
 mends a particular course of actions. In particular, we 524  
 consider what percentage of the input appears in the minimal 525  
 explanation. The smaller the explanation, the easier it is for a 526  
 human to interpret. We also focus on the scalability of Algo- 527  
 rithm 1 as the size of the policy increases. For reproducibility, 528  
 the repository [omitted for anonymity] provides our algo- 529  
 rithm implementation, benchmarks used, learnt policies, and 530  
 the scripts to learn them and run the experiments. 531

### 6.1 Experimental Setup 532

**Hardware** All experiments were run on a machine with an 533  
 AMD Ryzen Threadripper 3990X CPU, with 64 cores/128 534  
 threads, a clock speed of 2.9 GHz base, 4.3 GHz max boost, 535  
 and 128 GB of memory. 536

**MIP Configuration** Gurobi is the MIP solver used for 537  
 the experiments. To ensure the model is accurate enough 538  
 for our experiments, we set the integer feasibility tolerance 539  
 (`IntFeasTol`) to  $10^{-9}$  and the error for function approx- 540  
 imations (`FuncPieceError`) to  $10^{-6}$ . Presolve is also 541  
 turned off, as our use case is assessing the feasibility of 542  
 the model rather than finding an optimal value. If a single 543  
 call to MIP exceeded 2 hours, the algorithm was terminated 544  
 for the particular problem instance. 545

**ASNet policies** We took all deterministic domains and 546  
 training instances from the code distributions of (Toyer et al. 547  
 2020) and (Steinmetz et al. 2022). We ran the script provided 548  
 by Toyer using 1 core, an 8h timeout, and 64gb of memory 549  
 to learn, from these instances, two-layer (i.e. two proposi- 550  
 tion layers and three action layers) *sparse* policies with skip 551  
 connections and without heuristic inputs. As described in 552  
 (Toyer et al. 2020, Sec. 6), the  $\ell_1$ -regulariser used to train 553  
 sparse policies results in many modules having coefficient 554  
 so close to zero that they are insignificant to the output of 555  
 the network and are pruned by the ASNet sparsification 556  
 procedure. Hence, these modules do not appear in the MIP 557  
 model, making a simpler model to solve. 558

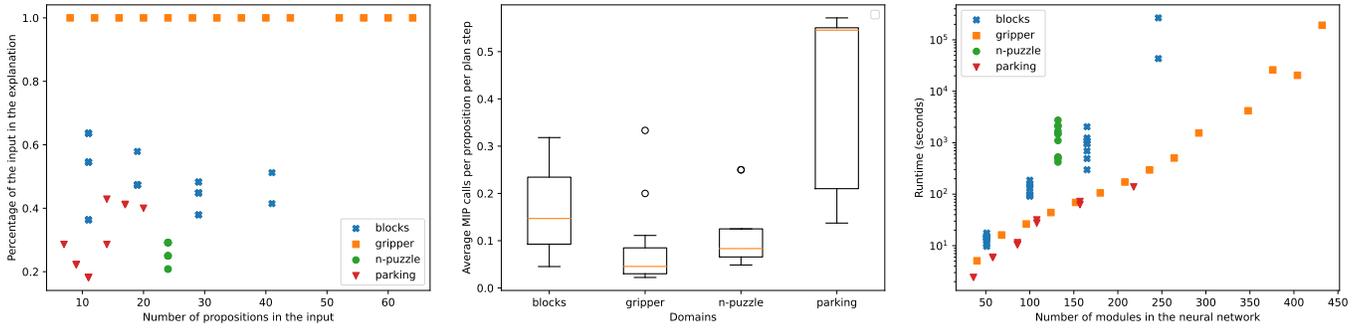


Figure 3: Experimental Results

## 6.2 Domains and Problems

We kept the domains for which the learnt policy could solve any problem within 150 execution steps, amongst a test set of 20 randomly generated small problems for that domain. This resulted in 4 policies for the domains of Blocksworld, Gripper,  $n$ -Puzzle, and Parking, which are present in the ASNet distribution. In order to evaluate our approach, we randomly generated, for each of these domains, the set of problems described below. Explanations were only computed for problems where the policy was able to reach the goal.

**Blocksworld** Experiments were conducted on 10 problems of each size (2, 3, 4 and 5 blocks) for a total of 40 problems. The policy was able to solve 90% of problems, and out of these, 83% had explanations computed within the time limit.

**Gripper** Explanations were computed using problems with 2 rooms and 1-15 balls. The policy was able to solve all problems, and explanations could be computed for 93%.

**$n$ -Puzzle** Experiments were conducted on all solvable and non-trivial combinations of 3-puzzle problems (on a 2x2 grid). Out of these 11 problems, the policy was correctly able to solve all instances and compute all explanations.

**Parking** Problems in the parking domain comprised of 2 cars and 2-4 curbs. Out of the 22 total problems, the policy reached the goal for 55% and we were able to compute explanations for 92% of these within the time limit.

## 6.3 Results

Figure 3 shows the size of the explanation as a percentage of the input (left-hand side graph), the average number of calls to the MIP solver per proposition and step in the plan (middle graph), and the CPU time of the algorithm (on a log scale) as a function of the size of the neural network policy for the problem (right-hand side graph).

**Size of Explanations** In all domains except Gripper, Algorithm 1 produces an explanation smaller than the ASNet’s input. For  $n$ -Puzzle, its small explanations are due to many static propositions. The grid in the problem is defined as propositions listing the neighbours of each position. As this grid is static, this will not appear in the explanation. On the other hand, no non-trivial explanations were found for Gripper, as the position of the gripper, the location of the

balls, and whether or not the gripper is holding anything are all relevant to the plan. Propositions stating which ball the gripper is carrying are not necessary, but they have already been removed from the MIP model due to the sparsity of the network.

**Number of Consistency Tests.** The average number of calls to the MIP solver per proposition and per action in the plan is always between 0 and 1. It represents how long propositions survive through the inner loop before being reinstated in the explanation or definitely ruled out. When that number is low, the algorithm quickly decides to keep the propositions in the explanation because of satisfiable consistency tests, or quickly decides to remove them because they are achieved by the effect of one of the first few actions. Gripper falls in the former camp: propositions are quickly ruled out as needing to be kept by the consistency test.  $n$ -puzzle falls in the latter camp: propositions quickly get removed from the explanation. In Blocksworld, and especially Parking, many propositions are removed from the explanation, but this happens much later in the sequence.

**Scalability** Much of the runtime is spent in MIP consistency tests. As can be seen from the size of problems we can address, reasoning about neural networks using MIP is a serious bottleneck. Across domains, we have observed that, as input propositions progressively get removed from the explanation, consistency tests generally become harder. This is because proving inconsistency generally becomes harder, whilst consistency becomes easier, but proving consistency leads us to reinsert the proposition, so the algorithm never reaches the easy region of the consistent side of the phase boundary. This behaviour is common with optimisation problems.

There are differences in scalability amongst the domains however. Gripper scaled the best, which is due to each iteration of the inner loop terminating early, as can be seen from the right-hand graph of Figure 3. The many MIP calls to the experiments for Parking problems increased the runtime significantly, resulting in the timeout of larger problems.

## 7 Conclusion, Limits, and Future Work

We have extended abductive explanations to sequences of decisions recommended by neural network policies. Our de-

640 composition approach to find a single minimal explanation  
641 incurs no overhead in comparison with the single decision  
642 case, once the length of the sequence is factored in.

643 Our approach makes a number of limiting assumptions  
644 which we discuss here together with possible extensions and  
645 future work. The first assumption is the availability of a  
646 symbolic planning model. An interesting avenue for future  
647 work is the extension of this approach to **learnt planning**  
648 **models**, also represented as neural networks.

649 We also assumed that actions have simple, unconditional  
650 effects. **Conditional effects** can be handled by our naive algo-  
651 rithm. However, Theorem 1 does not allow for them because  
652 it is impossible to apply conditional effects to a partial state.  
653 We assumed that we have no **background knowledge**, i.e.  
654 constraints that restrict the set of possible states (Thiébaux,  
655 Hoffmann, and Nebel 2005; Rintanen 2017). These can  
656 greatly simplify explanations. Yu et al. (2023) showed that in  
657 the single decision case, adding the background knowledge  
658  $K$  as another conjunct to the consistency tests performed  
659 by the greedy algorithm preserves the minimality of expla-  
660 nations. This property carries over to the multiple decision  
661 case and our naive algorithm. However, our decomposition  
662 algorithm may not return a minimal explanation with this aug-  
663 mented consistency test, because the set of reachable states  
664 cannot be represented by intersections of partial states with  
665  $K$ . We leave an efficient treatment of conditional effects and  
666 background knowledge for future work.

667 We have assumed that the actions and the policy are deter-  
668 ministic. Handling **stochastic actions and/or policies** could  
669 be achieved by generating explanations that pertain to a fi-  
670 nite tree of actions reachable with non-zero probability from  
671 the initial state. This would require applying progression  
672 and consistency tests at each branch of the tree. Handling  
673 stochastic policies would additionally require a more com-  
674 plex consistency test, as the decision  $d$  becomes a probability  
675 distribution over actions and the test must establish that it is  
676 not possible for the policy to return a different distribution.

677 Another avenue for future work is the generation of all  
678 (set-inclusion) **minimal explanations** and of **minimum car-**  
679 **dinality explanations**. This is likely to require a different  
680 decomposition of the problem than the one presented here,  
681 as well as effective heuristics to guide search. Finally, even  
682 in the single decision case, methods for computing formal  
683 explanations of neural networks suffer from scalability issues  
684 due to the expensive consistency test. New breakthroughs in  
685 MIP and SMT methods for analysing neural networks, new  
686 problem abstractions, and approximate explanation methods  
687 will be needed for these approaches to flourish.

## 688 References

689 Bäckström, C.; and Nebel, B. 1995. Complexity Results for  
690 SAS+ Planning. *Comput. Intell.*, 11: 625–656.  
691 Chakraborti, T.; Sreedharan, S.; and Kambhampati, S. 2020.  
692 The Emerging Landscape of Explainable Automated Plan-  
693 ning & Decision Making. In *Proc. IJCAI*, 4803–4811.  
694 Chakraborti, T.; Sreedharan, S.; Zhang, Y.; and Kambham-  
695 pati, S. 2017. Plan Explanations as Model Reconciliation:  
696 Moving Beyond Explanation as Soliloquy. In *Proc. IJCAI*,  
697 156–163.

Darwiche, A.; and Hirth, A. 2020. On the Reasons Behind 698  
Decisions. In *Proc. ECAI*, 712–720. 699  
Eifler, R.; Steinmetz, M.; Torralba, Á.; and Hoffmann, J. 700  
2020. Plan-Space Explanation via Plan-Property Depend- 701  
encies: Faster Algorithms & More Powerful Properties. In *Proc.* 702  
*IJCAI*, 4091–4097. 703  
Ferber, P.; Geißer, F.; Trevizan, F. W.; Helmert, M.; and 704  
Hoffmann, J. 2022. Neural Network Heuristic Functions for 705  
Classical Planning: Bootstrapping and Comparison to Other 706  
Methods. In *Proc. ICAPS*, 583–587. 707  
Ferber, P.; Helmert, M.; and Hoffmann, J. 2020. Neural 708  
Network Heuristics for Classical Planning: A Study of Hy- 709  
perparameter Space. In *Proc. ECAI*, 2346–2353. 710  
Fischetti, M.; and Jo, J. 2018. Deep neural networks and 711  
mixed integer linear optimization. *Constraints*, 23(3): 296– 712  
309. 713  
Garg, S.; Bajpai, A.; and Mausam. 2020. Symbolic Network: 714  
Generalized Neural Policies for Relational MDPs. In *Proc.* 715  
*ICML*, 3397–3407. 716  
Gehring, C.; Asai, M.; Chitnis, R.; Silver, T.; Kaelbling, L. P.; 717  
Sohrabi, S.; and Katz, M. 2022. Reinforcement Learning 718  
for Classical Planning: Viewing Heuristics as Dense Reward 719  
Generators. In *Proc. ICAPS*, 588–596. 720  
Groshev, E.; Goldstein, M.; Tamar, A.; Srivastava, S.; and 721  
Abbeel, P. 2018. Learning Generalized Reactive Policies 722  
Using Deep Neural Networks. In *Proc. ICAPS*, 408–416. 723  
Gurobi Optimization, LLC. 2023. Gurobi Optimizer Refer- 724  
ence Manual. [https://www.gurobi.com/wp-content/plugins/](https://www.gurobi.com/wp-content/plugins/hd_documentations/documentation/current/refman.pdf) 725  
[hd\\_documentations/documentation/current/refman.pdf](https://www.gurobi.com/wp-content/plugins/hd_documentations/documentation/current/refman.pdf). 726  
Ignatiev, A.; Narodytska, N.; and Marques-Silva, J. 2019. 727  
Abduction-Based Explanations for Machine Learning Mod- 728  
els. In *Proc. AAAI*, 1511–1519. 729  
Karia, R.; and Srivastava, S. 2021. Learning Generalized 730  
Relational Heuristic Networks for Model-Agnostic Planning. 731  
In *Proc. AAAI*, 8064–8073. 732  
Karia, R.; and Srivastava, S. 2022. Relational Abstractions for 733  
Generalized Reinforcement Learning on Symbolic Problems. 734  
In *Proc. IJCAI*, 3135–3142. 735  
Kasenberg, D.; Thielstrom, R.; and Scheutz, M. 2020. Gen- 736  
erating Explanations for Temporal Logic Planner Decisions. 737  
In *Proc. ICAPS*, 449–458. 738  
Krarup, B.; Krivic, S.; Magazzeni, D.; Long, D.; Cashmore, 739  
M.; and Smith, D. E. 2021. Contrastive Explanations of 740  
Plans through Model Restrictions. *J. Artif. Intell. Res.*, 72:  
533–612. 741  
Lundberg, S. M.; and Lee, S. 2017. A Unified Approach 742  
to Interpreting Model Predictions. In *Proc. NeurIPS*, 4765–  
4774. 743  
Marques-Silva, J.; and Ignatiev, A. 2022. Delivering Trust- 744  
worthy AI through Formal XAI. In *Proc. AAAI*, 12342–  
12350. 745  
Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. “Why 746  
Should I Trust You?”: Explaining the Predictions of Any 747  
Classifier. In *Proc. KDD*, 1135–1144. 748  
Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2018. Anchors: 749  
High-Precision Model-Agnostic Explanations. In *Proc. AAAI*,  
1527–1535. 750  
751  
752  
753  
754

755 Rintanen, J. 2017. Schematic Invariants by Reduction to  
756 Ground Invariants. In *Proc. AAAI*, 3644–3650.

757 Shen, W.; Trevizan, F. W.; and Thiébaux, S. 2020. Learning  
758 Domain-Independent Planning Heuristics with Hypergraph  
759 Networks. In *Proc. ICAPS*, 574–584.

760 Sreedharan, S.; Chakraborti, T.; and Kambhampati, S. 2021.  
761 Foundations of explanations as model reconciliation. *Artif.*  
762 *Intell.*, 301: 103558.

763 Steinmetz, M.; Fiser, D.; Eniser, H. F.; Ferber, P.; Gros, T. P.;  
764 Heim, P.; Höller, D.; Schuler, X.; Wüstholtz, V.; Christakis,  
765 M.; and Hoffmann, J. 2022. Debugging a Policy: Automatic  
766 Action-Policy Testing in AI Planning. In *Proc. ICAPS*, 353–  
767 361.

768 Thiébaux, S.; Hoffmann, J.; and Nebel, B. 2005. In defense  
769 of PDDL axioms. *Artif. Intell.*, 168(1-2): 38–69.

770 Toyer, S.; Thiébaux, S.; Trevizan, F. W.; and Xie, L. 2020.  
771 ASNets: Deep Learning for Generalised Planning. *J. Artif.*  
772 *Intell. Res.*, 68: 1–68.

773 Toyer, S.; Trevizan, F. W.; Thiébaux, S.; and Xie, L. 2018.  
774 Action Schema Networks: Generalised Policies With Deep  
775 Learning. In *Proc. AAAI*, 6294–6301.

776 Vasileiou, S. L.; Yeoh, W.; Son, T. C.; Kumar, A.; Cashmore,  
777 M.; and Magazzeni, D. 2022. A Logic-Based Explanation  
778 Generation Framework for Classical and Hybrid Planning  
779 Problems. *J. Artif. Intell. Res.*, 73: 1473–1534.

780 Vinzent, M.; Steinmetz, M.; and Hoffmann, J. 2022. Neural  
781 Network Action Policy Verification via Predicate Abstraction.  
782 In *Proc. ICAPS*, 371–379.

783 Younes, H.; and Littman, M. 2004. PPDDL1.0: An extension  
784 to PDDL for expressing planning domains with probabilistic  
785 effects. Technical report, CMU.

786 Yu, J.; Ignatiev, A.; Stuckey, P. J.; Narodytska, N.; and  
787 Marques-Silva, J. 2023. Eliminating the Impossible, What-  
788 ever Remains Must be True. In *Proc. AAAI*.

789 Zhang, Z.; and Geißer, F. 2022. Extending Graph Neural Net-  
790 works for Generalized Stochastic Planning. In *Proc. ICAPS-  
791 21 Planning and Reinforcement Learning Workshop (PRL)*.

## 792 A Proof of Theorem 1

793 Before proving Theorem 1, we need one additional definition.  
794 Given a partial state  $s$  and an action  $a$  such that for all  $x \in$   
795  $X_s \cap X_{\text{eff}(a)}$ ,  $s[x] = \text{eff}(a)[x]$ , we define the *regression* of  
796  $s$  through  $a$  as the smallest partial state in which applying  
797 action  $a$  leads to a state satisfying  $s$ :

$$\text{reg}_a(s) = (s \ominus \text{eff}(a)) \oplus \text{pre}(a)$$

798 **Lemma 1.** *If  $\mathbf{z}$  explains  $a_1, \dots, a_n$ , then for all  $i \in$*   
799  *$\{0, \dots, n-1\}$ ,  $\mathbf{z}_{i+1} := \text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z})$  explains  $a_{i+1}$ .*

800 We shall prove that for all  $i \in \{0, \dots, n-1\}$ , for all state  
801  $s_{i+1}$  that completes  $\mathbf{z}_{i+1}$ , there exists  $s_i$  such that i)  $a_i$  is  
802 applicable in  $s_i$ , ii)  $s_i \xrightarrow{a_i} s_{i+1}$  is a transition of the planning  
803 domain, and iii)  $s_i$  completes  $\mathbf{z}_i$ . By recursion, we end up  
804 with a state  $s_1$  that completes  $\mathbf{z}_1$  such that  $s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2}$   
805  $\dots \xrightarrow{a_i} s_{i+1}$ . Since  $s_1$  completes  $\mathbf{z}_1 = \mathbf{z}$ , we know that

$\pi^n(s_1) = a_1, \dots, a_n$ . Therefore,  $\pi(s_{i+1}) = a_{i+1}$ . As this is  
true for all  $s_{i+1}$ ,  $\mathbf{z}_{i+1}$  explains  $a_{i+1}$ .

The proof will be done by induction, i.e., we assume that it  
is true for  $i$ . (Base case for  $i = 0$  is trivially true as  $\mathbf{z}_1 = \mathbf{z}$ .)

Given our state  $s_{i+1}$ , we choose  $s_i$  as one of the states that  
complete  $\mathbf{z}_i \oplus \text{reg}_{a_i}(s_{i+1})$ . That is

$$s_i \supseteq \mathbf{z}_i \oplus ((s_{i+1} \ominus \text{eff}(a_i)) \oplus \text{pre}(a_i)). \quad (2)$$

We prove that  $s_i$  satisfies the three points above.

i) Eq. 2 clearly enforces the precondition  $\text{pre}(a_i)$ , so  $a_i$  is  
indeed applicable in  $s_i$ .

ii) We note

$$\text{prg}_{a_i}(s_i) \supseteq (s_{i+1} \ominus \text{eff}(a_i)) \oplus \text{pre}(a_i) \oplus \text{eff}(a_i).$$

which can be simplified by

$$\text{prg}_{a_i}(s_i) \supseteq s_{i+1} \oplus \text{pre}(a_i) \oplus \text{eff}(a_i). \quad (3)$$

We also note that all variables  $x \in X$  are specified in  
 $\text{prg}_{a_i}(s_i)$ , regardless of the choice of  $s_i$ . Consider any vari-  
able  $x$ ; does  $\text{prg}_{a_i}(s_i)[x] = s_{i+1}[x]$  hold?

- If  $x \in X_{\text{eff}(a_i)}$ , then  $\text{prg}_{a_i}(s_i)[x] = \text{eff}_{a_i}(x)$  and, since  
 $\mathbf{z}_{i+1} = \text{prg}_{a_i}(\mathbf{z}_i)$ ,  $s_{i+1}[x] = \mathbf{z}_{i+1}[x] = \text{eff}_{a_i}(x)$ .

- If  $x \in X_{\text{pre}(a_i)} \setminus X_{\text{eff}(a_i)}$ , then  $\text{prg}_{a_i}(s_i)[x] = s_i(x) =$   
 $\text{pre}(a_i)[x]$  since  $a_i$  does not modify  $x$  and  $a_i$  is applicable  
in  $s_i$ .

Furthermore, we note that  $s_{i+1}$  completes  $\text{prg}_{a_1}(\mathbf{z}_i)$  and  
that  $\mathbf{z}_i[x] = \text{pre}(a_i)[x]$  since (by induction)  $\mathbf{z}_i$  explains  
 $a_i$ ; therefore  $s_{i+1}[x] = \text{pre}(a_i)[x] = \text{prg}_{a_i}(s_i)[x]$ .

- If  $x \in X \setminus X_{\text{pre}(a_i)} \setminus X_{\text{eff}(a_i)}$ , then  $\text{prg}_{a_i}(s_i)[x] =$   
 $s_i[x] = s_{i+1}[x]$ .

So  $\text{prg}_{a_1}(s_i) = s_{i+1}$ .

iii) Does  $s_i$  complete  $\mathbf{z}_i$ ? Since, by construction,  $s_i$  completes  
 $\mathbf{z}_i \oplus ((s_{i+1} \ominus \text{eff}(a_i)) \oplus \text{pre}(a_i))$ , the question is whether  
some assignment in the right operand of  $\mathbf{z}_i \oplus$  contradicts an  
assignment in  $s_i$ . We know that  $\mathbf{z}_i$  explains  $a_i$ , so  $\text{pre}_{a_i}$  will  
not contradict  $\mathbf{z}_i$ . State  $s_{i+1}$  completes  $\mathbf{z}_{i+1}$  which differs  
with  $\mathbf{z}_i$  only on  $\text{eff}(a_i)$ . However, the expression above re-  
moves  $\text{eff}(a_i)$  from  $s_{i+1}$ , so that the right operand does not  
map any variable to a different value than  $\mathbf{z}_i$ .

**Lemma 2.** *Let  $\mathbf{z}$  be a partial state, and for all  $i \in \{0, \dots, n-$*   
840  *$1\}$ ,  $\mathbf{z}_{i+1} := \text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(\mathbf{z})$ . If for all  $i$ ,  $\mathbf{z}_i$  explains  $a_i$ ,*  
841 *then  $\mathbf{z}$  explains  $a_1, \dots, a_n$ .*

Assume that  $\mathbf{z}_i$  explains  $a_i$  for all  $i$ . Let  $s$  be a state com-  
pleting  $\mathbf{z}$  and let  $s_1 \xrightarrow{a_1} \dots \xrightarrow{a_n} s_n$  be the trajectory obtained  
by applying the sequence of actions  $a_1, \dots, a_n$  from  $s_1 = s$ .  
We shall prove  $\pi(s_i) = a_i$  for all  $i$  (which proves that  $\pi$   
recommends  $a_1, \dots, a_n$ ); this is proven by showing that  $s_i$   
completes  $\mathbf{z}_i$ .

The state  $s_{i+1}$  is defined as  $\text{prg}_{a_i} \circ \dots \circ \text{prg}_{a_1}(s)$ . Do we  
have  $\forall x \in X_{\mathbf{z}_{i+1}}$ ,  $s_{i+1}[x] = \mathbf{z}_{i+1}[x]$ ? Let  $j$  be the largest  
index in  $\{1, \dots, i\}$  such that  $x \in X_{\text{eff}(a_j)}$ . If  $j$  does not  
exist, then  $s_{i+1}[x] = s[x] = \mathbf{z}[x] = \mathbf{z}_{i+1}[x]$ . Otherwise,  
 $s_{i+1}[x] = \text{eff}(a_j)[x] = \mathbf{z}_{i+1}[x]$ . Either way, the variables  
of  $s_{i+1}$  map to the same value as those of  $\mathbf{z}_{i+1}$ . Therefore,  
 $\mathbf{z}_{i+1}$  explains  $a_{i+1}$  in  $s_{i+1}$ .

Theorem 1 is a consequence of Lemmas 1 and Lemma 2.