

# A New Approach to Plan-Space Explanation: Analyzing Plan-Property Dependencies in Oversubscription Planning

Rebecca Eifler,<sup>1</sup> Michael Cashmore,<sup>2</sup> Jörg Hoffmann,<sup>1</sup>  
Daniele Magazzeni,<sup>3</sup> Marcel Steinmetz<sup>1</sup>

<sup>1</sup>Saarland University, Saarland Informatics Campus, Saarbrücken, Germany,

<sup>2</sup>University of Strathclyde, Computer and Information Sciences, Glasgow, UK,

<sup>3</sup>King’s College London, Department of Informatics, London, UK,

{lastname}@cs.uni-saarland.de, {firstname.lastname}@strath.ac.uk, {firstname.lastname}@kcl.ac.uk

## Abstract

In many usage scenarios of AI Planning technology, users will want not just a plan  $\pi$  but an explanation of the space of possible plans, justifying  $\pi$ . In particular, in oversubscription planning where not all goals can be achieved, users may ask why a conjunction  $A$  of goals is not achieved by  $\pi$ . We propose to answer this kind of question with the goal conjunctions  $B$  excluded by  $A$ , i. e., that could not be achieved if  $A$  were to be enforced. We formalize this approach in terms of plan-property dependencies, where plan properties are propositional formulas over the goals achieved by a plan, and dependencies are entailment relations in plan space. We focus on entailment relations of the form  $\bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$ , and devise analysis techniques *globally* identifying all such relations, or *locally* identifying the implications of a single given plan property (user question)  $\bigwedge_{g \in A} g$ . We show how, via compilation, one can analyze dependencies between a richer form of plan properties, specifying formulas over action subsets touched by the plan. We run comprehensive experiments on adapted IPC benchmarks, and find that the suggested analyses are reasonably feasible at the global level, and become significantly more effective at the local level.

## 1 Introduction

AI plan generation technology serves to generate a plan  $\pi$  based on a world model. If the model and optimization objective reflect the real world and user preference with absolute accuracy, then  $\pi$  can be executed as-is. Yet in many usage scenarios, that is not so. Often, models are approximate, optimization objectives are highly complex and/or implicit in the heads of the human users, and bad plan decisions can be highly detrimental. A prominent example are space applications as discussed e. g. by Smith (2012), but arguably also many other applications ranging from Industry 4.0 to robot-aided disaster recovery. As Smith pointed out, planning should then be an iterative process in which the planning system suggests plan candidates, to be inspected and criticised by human users for iterative plan improvement.

As Smith also pointed out, plan explanation is a crucial step of such an iterative process, as a key means for user inspection. In particular, questions of the form “Why does the plan not achieve goal  $G$ ?” or “Why don’t you satisfy

preference  $P$ ?” need to be answered. Such answers require insights about the space of possible plans. We propose to address this by **plan-property entailments**, as in “Because achieving  $G$  would necessitate to either forego  $G'$  or use  $> 100$  energy units”. More formally, a plan property  $p$  is a Boolean function on plans (e. g.  $G$  true at end), and  $p$  entails  $q$  if all plans that satisfy  $p$  also satisfy  $q$ . A user question of the form “Why not  $p$ ?” can then be answered in terms of those  $q$  entailed by  $p$  (e. g.  $G'$  false at end or energy usage  $> 100$ ). We will refer to this as a **local explanation**, whereas a **global explanation** shows the graph of all plan-property entailments. The former makes sense in iterative planning as outlined, the latter makes sense when an overview of plan space is desirable. (For example, in simulated penetration testing (Boddy et al. 2005; Hoffmann 2015), an overview of the space of possible attacks may be useful.)

Our form of explanation falls into the class of *contrastive* explanations as discussed e. g. by Miller (2019). Previous work on local explanation, i. e., answering questions of the form “Why not  $p$ ?”, has suggested to generate a new plan  $\pi'$  that satisfies  $p$ , and answer the question based on comparing  $\pi$  and  $\pi'$  (Smith 2012; Fox, Long, and Magazzeni 2017). A weakness of this idea is that there may be differences between  $\pi$  and  $\pi'$  unrelated to  $p$ . Our approach replaces the *existential* answer generating a single alternative plan  $\pi'$  with a *universal* answer determining shared properties of *all* possible such alternatives. From this perspective, our proposal is a new, universal variant of contrastive plan explanation.

The concept of plan-property entailments is generic and in principle applicable to many AI Planning contexts. Here, we instantiate it for goal dependencies in oversubscription planning (OSP, e. g. (Smith 2004; Domshlak and Mirkis 2015)), where not all goals can be achieved. The plan properties we consider are goal-fact conjunctions  $\bigwedge_{g \in A} g$  or negations  $\neg \bigwedge_{g \in B} g$  thereof, and we identify **goal exclusion** entailments  $\bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$  stating that, if a plan achieves all goals in  $A$ , then it must forego at least one goal in  $B$ .

We spell out our framework for this form of plan properties and entailment relations. We introduce algorithms for computing local and global explanations, leveraging and extending recent nogood learning methods (Steinmetz and Hoffmann 2017b; 2017a). We show that the same framework can address a richer form of plan properties, namely **action-**

**set properties** defined as propositional formulas over atoms  $\{A_1, \dots, A_n\}$  asking whether an action subset  $A_i$  is touched by the plan. Action-set properties can be compiled into goal facts relatively easily, leading to an effective entailment analysis. We run comprehensive experiments on IPC benchmarks adapted to OSP as in previous work (Domshlak and Mirkis 2015; Katz et al. 2019), and on a collection of benchmarks we extended with action-set properties. We find that global explanations are reasonably feasible to compute, compared to OSP and to optimal classical planning. Computing local explanations is significantly easier.

Some related work will be discussed near the end of the paper (Section 7), to not interrupt the text flow.

## 2 Background

We consider the *finite-domain representation (FDR)* framework (Bäckström and Nebel 1995; Helmert 2009). An FDR task  $\tau$  is a tuple  $\tau = (V, A, c, I, G)$  where  $V$  is the set of **variables**,  $A$  is the set of **actions**,  $c : A \rightarrow \mathbb{R}_0^+$  is the action **cost** function,  $I$  is the **initial state**, and  $G$  is the **goal**. A **state**, in particular  $I$ , is a complete assignment to  $V$ ;  $G$  is a partial assignment to  $V$ ; each action  $a \in A$  has a **precondition**  $pre_a$  and an **effect**  $eff_a$ , both partial assignments to  $V$ . We will refer to variable-value pairs  $v = d$  as **facts**, and we will identify partial variable assignments with sets of facts. An action  $a$  is **applicable** in a state  $s$  if  $pre_a \subseteq s$ . The outcome state  $s[[a]]$  is like  $s$  except that  $s[[a]](v) = eff_a(v)$  for those  $v$  on which  $eff_a$  is defined. The outcome state of an iteratively applicable action sequence  $\pi$  is denoted by  $s[[\pi]]$ .

An **oversubscription planning (OSP) task** is a tuple  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  like an FDR task but with two goals – the **hard** goal  $G^{\text{hard}}$  and the **soft** goal  $G^{\text{soft}}$ , assumed to be defined on disjoint sets of variables – as well as an additional **cost bound**  $b \in \mathbb{R}_0^+$ . A **plan** is an action sequence  $\pi$  whose summed-up cost is  $\leq b$  and where  $G^{\text{hard}} \subseteq I[[\pi]]$ . Good plans achieve a “maximally valuable” subset of  $G^{\text{soft}}$ . The latter is usually defined based on goal-fact rewards. Here we assume instead that the user’s preferences over  $G^{\text{soft}}$  are difficult to specify and/or elicitate, and planning is an iterative process as described by Smith (2012). We are concerned with the analysis of plan-property dependencies within  $G^{\text{soft}}$  to support that process.

We will also consider planning with non-reproducible resources, encoded as discrete FDR state variables. This can be viewed as a special case of our OSP formulation. In our experiments, among others we will use resource-constrained planning (RCP) benchmarks by Nakhost et al. (2012), which allow to control resource constrainedness (basically the ratio between available vs. required resource amounts).

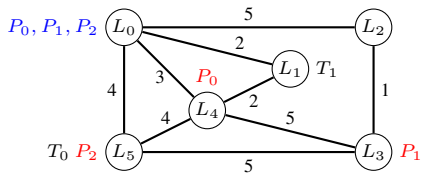


Figure 1: Running example from IPC NoMystery.

For illustration, we will use an RCP running example from IPC NoMystery, a transportation domain with fuel consumption. The variables encode truck and package positions, actions drive trucks or load/unload packages. Drive actions consume an amount of fuel (of the respective truck) depending on the road connection taken. As a concrete example, we will use the task with two trucks and three packages illustrated in Figure 1. Fuel consumptions are indicated at road segments. The packages are initially at  $L_0$  (shown in blue); their goal locations are  $L_4, L_3,$  and  $L_5$  (shown in red).

## 3 Goal-Property Dependencies in OSP

We next spell out our framework for plan properties, entailment relations between them, and the forms of explanations we aim at. We do so in FDR-based OSP as defined above, but in principle our definitions are generic and can be instantiated for arbitrary plan properties and planning frameworks. We will discuss the definitions from that perspective.

### 3.1 Plan Properties and Property Entailment

The plan properties we consider here are formulas over the soft goals:

**Definition 1 (Plan Properties).** Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.

A **plan property** is a function  $p_\phi : \Pi \rightarrow \{\text{true}, \text{false}\}$  where  $\phi$  is a propositional formula over the atoms  $G^{\text{soft}}$ , and  $p_\phi(\pi) = \text{true}$  iff  $\phi$  evaluates to true under the truth value assignment where  $g \in G^{\text{soft}}$  is true iff  $g \in I[[\pi]]$ .  $p_\phi$  is **conjunctive** if  $\phi$  has the form  $\bigwedge_{g \in A} g$  or  $\neg \bigwedge_{g \in B} g$ .

We identify  $p_\phi$  with the characterizing formula  $\phi$ . Our analyses are formulated for conjunctive plan properties (which through compilation can address arbitrary  $p_\phi$ ).

In general, a plan property can be any function mapping a task and an action sequence to a Boolean value. Examples are temporal plan trajectory constraints, deadlines, bounds on resource consumption, and the aforementioned action-set properties. To the extent that such properties can be compiled into goal facts, conjunctive plan properties can be used to analyze their dependencies. We will explore this possibility here with a compilation for action-set properties.

We formalize plan-property dependencies as entailments in the space of plans  $\Pi$ :

**Definition 2 ( $\Pi$ -Entailment).** Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.

We say that  $\pi \in \Pi$  **satisfies** a plan property  $\phi$ , written  $\pi \models \phi$ , if  $p_\phi(\pi) = \text{true}$ . We denote by  $\mathcal{M}_\Pi(\phi) := \{\pi \mid \pi \in \Pi, \pi \models \phi\}$  the subset of plans that satisfy  $\phi$ . We say that  $\phi$   **$\Pi$ -entails**  $\psi$ , written  $\Pi \models \phi \Rightarrow \psi$ , if  $\mathcal{M}_\Pi(\phi) \subseteq \mathcal{M}_\Pi(\psi)$ .

This definition views  $\Pi$  in the role traditionally taken by a knowledge base, identifying a set of “possible worlds” within which entailment over plan properties is considered. Observe that, given this,  $\Pi$ -entailment is more than standard entailment:  $\phi \Rightarrow \psi$  implies that  $\Pi \models \phi \Rightarrow \psi$ , but not vice versa.  $\Pi$ -entailment captures entailments specific to the space of plans  $\Pi$ . For example, in our illustrative NoMystery task, say that all goals are soft,  $T_0$  has initial fuel supply 13, and  $T_1$  has no fuel. Then  $\Pi \models at(P_0, L_4) \Rightarrow$

$\neg(at(P_1, L_3) \wedge at(P_2, L_5))$  because, if we achieve the goal for  $P_0$ , there is insufficient fuel to transport both other packages. If we set the initial fuel supply of  $T_0$  to 16, on the other hand, then the knowledge base changes –  $\Pi$  becomes more permissive – and that entailment no longer holds.

Note that the definition of  $\Pi$ -entailment is agnostic to the specification of  $\Pi$ . The definition applies unchanged to arbitrary planning frameworks and plan sets  $\Pi$ .

Our primary focus here will be on goal exclusions:

**Definition 3** (Goal Exclusions). *Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.*

A **goal exclusion** is an entailment of the form  $\Pi \models \bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$ . The entailment is **non-dominated** if there is no pair  $(A', B')$  where  $A' \subseteq A$ ,  $B' \subseteq B$ ,  $(A', B') \neq (A, B)$ , and  $\Pi \models \bigwedge_{g \in A'} g \Rightarrow \neg \bigwedge_{g \in B'} g$ . The entailment is **non-rhs-dominated** if there is no  $B'$  where  $B' \subsetneq B$  and  $\Pi \models \bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B'} g$ .

Goal exclusions are of interest in OSP as they reflect the detailed (soft-)goal trade-offs in the user’s quest for a good plan. A non-dominated goal exclusion has subset-minimal  $A$  and  $B$ . This dominates entailments with larger  $A$  and/or  $B$  as it has a weaker left-hand side  $A$  (smaller conjunction) entailing a stronger right-hand side  $\neg B$  (smaller disjunction, after moving the negation inside). If  $A$  is fixed, then only the right-hand side  $B$  needs be minimal. We will use non-dominated entailments to give more compact explanations.

### 3.2 Local and Global Explanations

As previously hinted, we propose to employ the concept of plan-property entailment for the purpose of giving local and global explanations of the plan space  $\Pi$ .

For local explanations, we assume a user question of the form “Why do you not achieve property  $\phi$ ?”, which we answer with the set of plan properties  $\psi$  entailed by  $\phi$ :

**Definition 4** (Local Explanation (LE)). *Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.*

For a plan property  $\phi$ , the **local explanation (LE)** for  $\phi$  is the set  $\{\psi \mid \Pi \models \phi \Rightarrow \psi\}$  of plan properties  $\Pi$ -entailed by  $\phi$ . For  $\phi = \bigwedge_{g \in A} g$ , the **goal-exclusion LE** for  $\phi$  is  $\{\psi \mid \psi = \neg \bigwedge_{g \in B} g, \Pi \models \phi \Rightarrow \psi \text{ is non-rhs-dominated}\}$ .

Such an answer makes sense if the entailed properties  $\psi$  are undesirable. This is the case, in particular, for goal-exclusion LEs. In our example with  $T_0$  fuel 13, the answer to  $\phi = at(P_0, L_4)$  “Why do you not achieve the goal for  $P_0$ ?” would be  $\psi = \neg(at(P_1, L_3) \wedge at(P_2, L_5))$  “Because that would necessitate to forego the goal for either  $P_1$  or  $P_2$ ”.

From a general perspective, plan properties here serve as an abstraction level at which to explain  $\Pi$  to a user. The underlying assumption is that  $\Pi$  is large and/or the mechanisms that generate  $\Pi$  are complex, so that an abstract form of explanation is needed. The abstraction level can be controlled through the number and granularity of plan properties. Given this, while here we simply talk about all formulas over soft-goal facts, it can make sense to instead fix a more specific set  $P$  of plan properties the user has a vested interest in (raising the new sub-problem how to choose  $P$ ).

Note that, if the user question “Why do you not achieve property  $\phi$ ?” refers to a concrete given plan candidate  $\pi$ , then a simpler variant of local explanation is to return only those entailed  $\psi$  where  $\pi \not\models \psi$ , i.e., currently false plan properties that would become true when enforcing  $\phi$ . This is easier to compute and yields smaller explanations. In our experiments, we will only use explanations as per Definition 4, yielding an upper bound on computational hardness and explanation size, and avoiding the bias incurred by a particular method for generating candidate plans  $\pi$ .

A canonical notion of global explanation arises directly from the above. Instead of showing the implications of one specific plan property, one can show all such implications:

**Definition 5** (Global Explanation (GE)). *Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.*

Denote by  $[\phi]_{\Pi} := \{\psi \mid \mathcal{M}_{\Pi}(\phi) = \mathcal{M}_{\Pi}(\psi)\}$  the  $\Pi$ -equivalence class of  $\phi$ . The **global explanation (GE)** for  $\tau$  is the strict partial order  $\Rightarrow_{\Pi}$  over the classes  $[\phi]_{\Pi}$  where  $[\phi]_{\Pi} \Rightarrow_{\Pi} [\psi]_{\Pi}$  iff  $[\phi]_{\Pi} \neq [\psi]_{\Pi}$  and  $\Pi \models \phi \Rightarrow \psi$ .

The only somewhat non-obvious design decision here is to group plan properties into equivalence classes. That said, while this definition makes sense at the formal level, it is of doubtful practical value. There can be many plan properties (here: all propositional formulas over  $G^{\text{soft}}$ ) and many equivalence classes (here: in the worst case, one class for every set of truth-value assignments to  $G^{\text{soft}}$ ). This makes it questionable whether the GE can be processed by a user.

It therefore makes sense to focus on more limited forms of GEs, and to find ways to represent these more compactly. Here, we focus on non-dominated goal exclusions:

**Definition 6** (Goal-Exclusion GE). *Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.*

The **goal-exclusion GE** for  $\tau$  is the strict partial order  $\Rightarrow_{\Pi}$  over conjunctive plan properties where  $\phi \Rightarrow_{\Pi} \psi$  iff  $\Pi \models \phi \Rightarrow \psi$  is a non-dominated goal exclusion.

Equivalence classes are not needed here because we consider only goal exclusions  $\Pi \models \bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$ , so plan properties cannot entail each other. Our results on benchmarks suggest that the goal-exclusion GE is reasonably feasible to compute, and is often small.

Consider again our example, with  $T_0$  fuel 13. Taking the freedom to notate goal facts by the respective packages,  $\Pi \models P_0 \Rightarrow \neg(P_1 \wedge P_2)$ ,  $\Pi \models P_1 \Rightarrow \neg(P_0 \wedge P_2)$ , and  $\Pi \models P_2 \Rightarrow \neg(P_0 \wedge P_1)$  are non-dominated goal exclusions. They are not the only ones though:  $\Pi \models P_0 \wedge P_1 \Rightarrow \neg P_2$  also satisfies the definition, and so does any way of distributing the three packages across the left-hand vs. right-hand sides of the entailment. So the size of the goal-exclusion GE, in terms of the ordering relations it contains, is eight.

Note that this GE is equivalent to the statement “ $G^{\text{soft}}$  is not solvable as a whole, but each of its subsets is”. Indeed, as we shall spell out next, for the special case of goal exclusions the GE is equivalent to the collection of minimal unsolvable goal subsets. We exploit this connection in our algorithms.

## 4 Computing Goal-Exclusion Explanations

We start with the algorithms for GEs, from which those for LEs follow easily.

## 4.1 Global Explanations

The special case of goal exclusions – though not the analysis of plan-property dependencies in general – boils down to the computation of **minimal unsolvable goal subsets (MUGS)**. A soft-goal set  $G \subseteq G^{\text{soft}}$  is a MUGS if  $G$  cannot be achieved but every  $G' \subsetneq G$  can:

**Proposition 1** (GE from MUGS). *Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task, and  $\Pi$  its set of plans.*

*Then  $\Pi \models \bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$ , and that entailment is non-dominated, if and only if  $A \cup B$  is a MUGS.*

*Proof.* A  $\Pi$ -entailment  $\Pi \models \bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$  clearly holds iff  $A \cup B$  is unsolvable. Non-dominated entailments result from set-inclusion minimal  $A$  and  $B$ , corresponding to the set-inclusion minimality of MUGS.  $\square$

We can hence represent goal-exclusion explanations more compactly, showing one MUGS  $G$  for each set of goal-exclusion entailments with  $A \cup B = G$ . In our example, the only MUGS is  $\{at(P_0, L_4), at(P_1, L_3), at(P_2, L_5)\}$ .

The question now is how to compute all MUGS. This can be done through a search over goal sets, that we refer to as **systematic weakening (SysW)**:

- (1) the start node of the search is  $G^{\text{soft}}$ ;
- (2) each search step selects an open node  $G$ , calls a planner to test whether  $G$  is solvable in  $\tau$ , caches the result, and expands  $G$  if it is unsolvable;
- (3) the children of a node  $G$  are those  $G' \subset G$  where  $|G'| = |G| - 1$ .

Upon termination, the MUGS are those nodes  $G$  all of whose children are solvable.

Dually, **systematic strengthening (SysS)** starts from  $\emptyset$ , with search steps expanding solvable nodes, and children adding one more goal fact. Upon termination, the MUGS can be easily obtained from the unsolvable search nodes.

In both algorithms, the worst-case number of planner calls is exponential in  $|G^{\text{soft}}|$ , i. e., the number of goals whose dependencies are being analyzed. Intuitively, SysW is better suited if solvable goal subsets are large (and thus are encountered early), while SysS is better suited if they are small (unsolvable subsets are encountered early). For illustration, in the example SysW expands the single search node  $G^{\text{soft}}$ , while SysS expands all solvable subsets before reaching  $G^{\text{soft}}$ .

Search enhancements are important for scalability. In particular, goal sets can be reached from the start node by permutations of goal-fact removal/addition steps. For effective results re-use, we assign goal sets unique integer IDs and fix an ascending expansion order. When expanding a node  $G$  and considering a possible child  $G'$ , there are three possibilities: 1.  $id(G') > id(G)$ :  $G'$  not seen yet, insert into open list; 2.  $id(G') < id(G)$  and result for  $G'$  cached: use that result; 3.  $id(G') < id(G)$  and result for  $G'$  not cached: then the search was previously cut off above  $G'$ , so we know that  $G'$  is solvable (SysW) respectively unsolvable (SysS).

We created synergy with recent nogood learning techniques, **conjunction learning** (Steinmetz and Hoffmann 2017b) and **trap learning** (Steinmetz and Hoffmann 2017a).

These techniques refine dead-end detection methods (nogoods) based on the unsolvable states encountered in state space search on a planning task. Trap learning identifies unescapable non-goal regions of the state space. Conjunction learning enriches the critical-path heuristic  $h^C$  with more atomic conjunctions, and learns clauses by dead-end-state reduction. As the children tasks in our searches are closely related to their parents, the refined nogoods are likely to be useful still. So we **transfer** the nogoods along search paths, resulting in iteratively stronger and stronger nogoods.

For both forms of learning, the nogoods learned depend on the goal, so that only some of the nogoods remain valid for transfer in SysW where children remove goals. To identify that nogood subset, in conjunction learning we remember the subset  $G$  of goals proved to be unreachable by  $h^C$  in a learning step. Clause validity is preserved so long as at least one member of  $G$  is still present. In trap learning, the validity test checks, for every node of the trap, whether the node is still either mutex with the goal, or is still proved unsolvable by the heuristic function.

We extended conjunction learning to deal with OSP plan-cost bounds, by checking not whether  $h^C(s) = \infty$  but whether  $h^C(s)$  exceeds the remaining available cost. For trap learning, no straightforward extension exists, as trap nodes represent sets of states and thus would need to reason about bounds on remaining available cost. As trap learning is not the focus of our work, we leave this open for future research and use trap learning on RCP benchmarks only.

We acknowledge that the basic structure of SysW and SysS is simple, and relates to various prior works addressing conflict analysis in one or the other form. Our contribution here consists in assembling the algorithms in a form suited to our purposes, plus the described enhancements. We briefly discuss related work in Section 7.

## 4.2 Local Explanations

For local explanation, we are given a soft-goal subset  $A$  whose non-rhs-dominated entailments we wish to identify. This corresponds to fixing  $A$  and minimizing only  $B$  – which in turn corresponds to moving  $A$  into the hard goals and computing MUGS on the modified task.

Precisely, assume an OSP task  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$ , with set of plans  $\Pi$ . Consider the modified task  $\tau' := (V, A, c, I, G^{\text{hard}} \cup A, G^{\text{soft}} \setminus A, b)$ . Clearly, we have the goal exclusion  $\Pi \models \bigwedge_{g \in A} g \Rightarrow \neg \bigwedge_{g \in B} g$  iff  $B$  is unsolvable in  $\tau'$ ; and  $B$  is minimal in  $\tau'$  iff the goal exclusion is non-rhs-dominated. So the MUGS in  $\tau'$  yield exactly the goal-exclusion LE for  $A$  as per Definition 4.

Observe that  $\tau'$  has fewer soft goals, so MUGS computation becomes easier, resulting in an advantage of local analysis over global analysis. That advantage grows with  $|A|$ , i. e., intuitively, with how specific the user question is. At the extreme end,  $A = G^{\text{soft}}$  and both SysW and SysS on  $\tau'$  contain a single search node testing solvability of  $G^{\text{hard}} \cup G^{\text{soft}}$ .<sup>1</sup>

<sup>1</sup>As previously indicated, given a concrete candidate plan  $\pi$  one could choose to restrict the answer to entailments  $\neg \bigwedge_{g \in B} g$  where  $\pi \models \bigwedge_{g \in B} g$ . This corresponds to removing all  $g$  with  $\pi \not\models g$  from  $G^{\text{soft}}$  in  $\tau'$ , further simplifying the computation.

## 5 Action-Set Properties

The analyses just described can, in principle, be used to analyze dependencies between arbitrary plan properties, so long as these can be compiled into goal facts. Given the well-known power of compilation in planning languages (e.g. (Gazen and Knoblock 1997; Nebel 2000; Edelkamp 2006; Palacios and Geffner 2009; Baier, Bacchus, and McIlraith 2009)), there is large potential in this idea. As an example, here we consider what we refer to as action-set properties:

**Definition 7** (Action-Set Properties). *Let  $\tau = (V, A, c, I, G^{\text{hard}}, G^{\text{soft}}, b)$  be an OSP task,  $\Pi$  its set of plans, and  $A_1, \dots, A_n \subseteq A$ .*

*An **action-set property** for  $\tau$  and  $A_1, \dots, A_n$  is a function  $p_\phi : \Pi \rightarrow \{\text{true}, \text{false}\}$  where  $\phi$  is a propositional formula over the atoms  $A_1, \dots, A_n$ , and  $p_\phi(\pi) = \text{true}$  iff  $\phi$  evaluates to true under the truth value assignment where  $A_i$  is true iff  $\pi$  contains at least one action from  $A_i$ .*

As before, we identify action-set properties  $p_\phi$  with the characterizing formulas  $\phi$ . Arguably, action-set properties are practically relevant. They allow to express things like “objective  $x$  is covered by satellite  $y$ ”, “route  $x$  is not used”, “passengers  $x$  and  $y$  ride in the same vehicle”, if these are desirable but could be traded against other soft preferences. At the same time, the simple syntax of action-set properties lends itself to effective compilation, as follows.

Given  $\tau$ ,  $\Pi$ , and  $A_1, \dots, A_n$  as in Definition 7, to obtain a compiled task  $\tau'$

- introduce Boolean flags  $isUsed_i$  that are initially *false* and set to *true* by any action from  $A_i$ ;
- introduce formula-evaluation state variables and actions evaluating each  $p_\phi$  based on the  $isUsed_i$  flags (following (Gazen and Knoblock 1997; Nebel 2000)), setting Boolean flags  $isTrue_\phi$  storing the outcome values;
- introduce a separate 1. *planning phase* vs. 2. *formula-evaluation phase*, and a switch action allowing to go from 1. to 2. if  $G^{\text{hard}}$  is satisfied.

Then the planning-phase prefixes in  $\tau'$  are in one-to-one correspondence with  $\Pi$ , and given such a prefix  $\pi$  the evaluation phase in  $\tau'$  can achieve  $isTrue_\phi$  iff  $p_\phi(\pi) = \text{true}$ .

Now say that we want to analyze the dependencies across a given set  $P$  of action-set properties (e.g. possible undesirable consequences of not using route  $X$ ). We are given  $\tau$ ,  $\Pi$ , and  $P$ ; we want to identify dependencies of the form  $\Pi \models \bigwedge_{\phi \in A} \phi \Rightarrow \neg \bigwedge_{\psi \in B} \psi$ . With the above, this can be done by considering  $\tau'$  with soft goals  $\{isTrue_\phi \mid \phi \in P\}$ , and identifying each  $isTrue_\phi$  with  $\phi$  in the outcome.

In the NoMystery domain action-set properties that can make sense are, for example, *uses  $T_i$  ( $L_x, L_y$ )* (truck  $T_i$  drives at least once from  $L_x$  to  $L_y$  or vice versa), *doesn't use  $T_i$  ( $L_x, L_y$ )* (the opposite), or *same truck  $P_x P_y$*  (both packages are delivered by the same truck). In our example task, say we fix the package destinations as hard goals, and we want to analyze the properties 1. uses  $T_0$  ( $L_2, L_3$ ); 2. same truck  $P_1 P_2$ ; 3. uses  $T_0$  ( $L_4, L_3$ ); 4. same truck  $P_2 P_0$ ; 5. doesn't use  $T_0$  ( $L_0, L_5$ ); 6. uses  $T_1$  ( $L_5, L_4$ ). Say initial fuel is 16 for  $T_0$  and 7 for  $T_1$ . Computing the MUGS over the soft

goals representing the six action-set properties, it turns out there are seven minimal unsolvable subsets of these properties, each of size three. A user could, for example, ask “Why do you not avoid the road  $L_0 - L_5$  (which has a lot of traffic at the moment)?”, translating into the question “Why do you not achieve property 5?”. One of the non-rhs-dominated entailments of property 5 is  $\neg(\text{property 2} \wedge \text{property 4})$ , corresponding to the answer “Because if you don't use that road, then you cannot deliver all packages with a single truck.”

We remark that, in a preliminary exploration, we implemented a compilation for LTL plan properties based on previous work (Edelkamp 2006; Baier, Bacchus, and McIlraith 2009). Our results are promising, yet indicate that algorithmic optimizations are needed to obtain good performance for complex properties. This remains a topic for future work.

## 6 Experiments

We implemented our approach in Fast Downward (FD) (Helmert 2006). We evaluate it, in turn, on IPC benchmarks modified for OSP planning, and on a selection of IPC benchmarks we extended with action-set properties.

The base planner called by our SysS and SysW algorithms on each search node runs forward search using  $h^{\text{FF}}$  (Hoffmann and Nebel 2001), optionally with conjunction or trap learning. The experiments were run on a cluster of Intel E5-2660 machines running at 2.20 GHz, with time (memory) cut-offs of 30 minutes (4 GB).

### 6.1 IPC-Based OSP Benchmarks

Following Katz et al. (2019), for every IPC benchmark task  $(V, A, c, I, G)$  with smallest known plan cost  $C$  as per `planning.domains` (Muisse 2016), we obtained three OSP tasks by setting the cost bound to  $b = x * C$  where  $x \in \{0.25, 0.5, 0.75\}$ . We used soft goals only, i.e.,  $G^{\text{soft}} = G$  and  $G^{\text{hard}} = \emptyset$ . For implementation reasons, we omitted tasks with  $\geq 32$  goals (which would be infeasible for MUGS analysis anyhow). We consider conjunction learning, but not trap learning as that cannot deal with OSP cost bounds.

**IPC Global Explanations** Figure 2 shows our data for global explanations, i.e., computing all MUGS. Note first the #MUGS data in the rightmost part of the table. As the data shows, the size of this global explanation is often small.

The rest of the figure focuses on computational performance. As a measure to compare against, we use optimal planning and OSP as reference points. The  $h^{\text{LM-cut}}$  column gives coverage for  $A^*$  with  $h^{\text{LM-cut}}$  (Helmert and Domshlak 2009) run on the original IPC instance without a cost bound, as a comparison to solvable optimal planning. The OSP column gives coverage for the most recent OSP planner (Katz et al. 2019). It is expected that our algorithms, solving a more complex problem, will perform worse than the reference points. The question is, *how much worse?*

As a short summary of the answer provided by Figure 2 to that question, compared to optimal planning, with the smallest cost bound  $x = 0.25$  our analysis is actually more effective. But that changes for larger cost bounds where the base planner's search space is larger and accordingly our analysis



take the image?” is answered with “Because then I can’t take any of the samples.”

- Woodworking instance 01(opt08): There are the goals 1.  $colour(p_0, mauve)$ , 2.  $colour(p_1, green)$ , 3.  $surface\_condition(p_0, verysmooth)$ , 4.  $surface\_condition(p_1, smooth)$ . Among the MUGS there are  $\{1, 4, 5\}$ ,  $\{2, 4, 5\}$ . So the question “Why don’t you grind  $p_0$  and  $p_1$ ?” is answered with “Because then I cannot paint them.”

**IPC Local Explanations** For local explanations – entailments of soft-goal conjunctions  $\bigwedge_{g \in A} g$  – we set up an experiment where we used only benchmark tasks with  $\geq 6$  goals (resulting in 37 non-empty domains), scaled question size  $|A|$  from 1 to 5, and randomly selected 10 questions of each size. Figure 3 shows the data for SysW with conjunction learning at the most challenging cost bound 0.75. Averages are taken over instances where all questions were solved, as otherwise the results are distorted by large instances solved only for high values of  $|A|$ . As expected, computational effort and explanation size decrease with  $|A|$ .

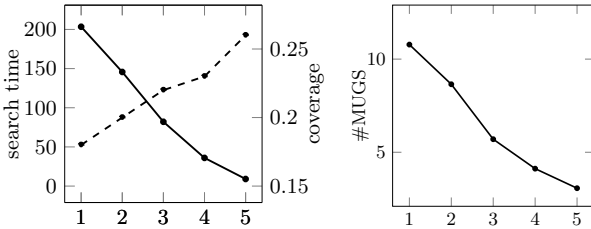


Figure 3: Local explanation results on IPC benchmarks, cost bound 0.75, for SysW as a function of  $|A|$ . Left: coverage, i. e., fraction of solved questions (dashed, right y-axis); and average runtime (solid, left y-axis). Right: average #MUGS.

These results are consistent across domains, and are much stronger in some cases. For example, while for  $|A| = 1$  our coverage is higher than that of  $h^{LM-cut}$  in only 2 domains here, for  $|A| = 5$  that is so in 6 domains. Compared against OSP, the number of domains with higher coverage rises from 5 to 13. Local explanations are typically small even for domains where global explanations can be large as shown in Figure 2: the average number of MUGS for  $|A| = 5$  is  $< 7$  in all but 3 domains, and is  $< 15$  except in Visitall (21.81).

## 6.2 Action-Set Properties

To evaluate the use of our framework with more complex plan properties, we experimented with the compilation of action-set properties as per Section 5. We extended four domains with action-set properties: NoMystery, Rovers, and TPP as per (Nakhost, Hoffmann, and Müller 2012), with controlled resource constrainedness (cf. Section 2); plus the Blocksworld as an intuitively differently structured domain. In NoMystery, the action-set properties are as in our illustrative example. In Rovers, the properties ask which rover or camera is used for which observation. In TPP, they ask which road segments are used, and which goods are bought at which markets. In Blocksworld, we include two hands and

the properties ask which hand is used for which blocks. We set the original goal facts as hard goals, so that the analysis determines exclusion relations over conjunctions of action-set properties.

We encoded resource consumption into the FDR state variables, enabling the use of trap learning which turns out to be highly beneficial here. We generated benchmarks of sizes around the feasibility borderline, and we experimented with resource constrainedness  $x \in \{1.0, 1.25, 1.5, 2.0\}$ . For reference, we ran A\* with  $h^{LM-cut}$  as before, on all goal facts (original plus compiled action-set properties). Current OSP planners cannot handle hard goals, so we used our base planner with trap learning instead as a second reference point.

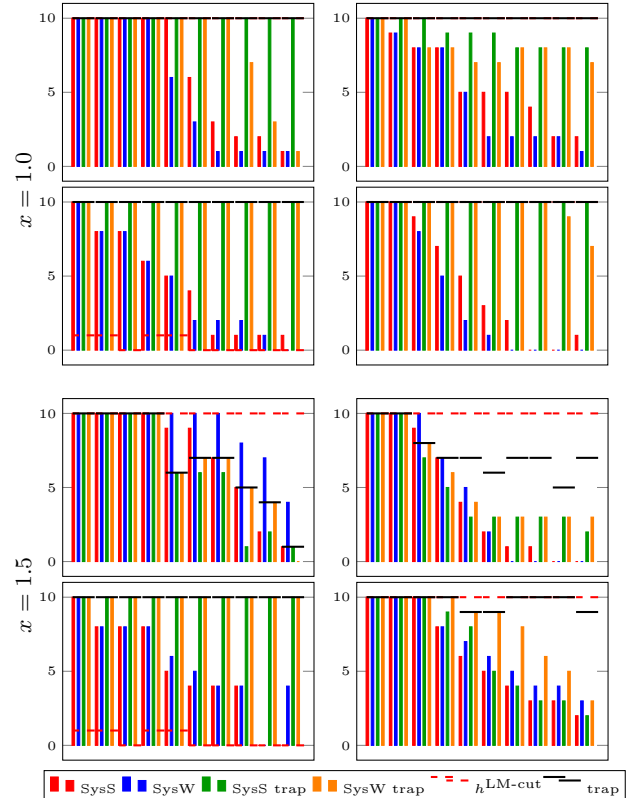


Figure 4: Coverage for global action-set property explanation at  $x = 1.0$  and  $x = 1.5$ . Number of properties on x-axis, scaling from 1 to 10. Blocksworld top left, NoMystery top right, Rovers bottom left, TPP bottom right.

Figure 4 shows performance data for global explanation. We consider  $x = 1.0$  and  $x = 1.5$  as representatives of the overall picture. At  $x = 1.0$ , computing global action-set property explanations is feasible compared to the reference points. SysS with trap learning matches reference point coverage except for moderate losses in NoMystery. As more resources becomes available, similarly to before our analysis problem becomes harder to solve, and the coverage gap to the reference points increases. Rovers is an exception to this because, there, all pairs of action-set properties are unsolvable, which trap learning discovers quite effectively.

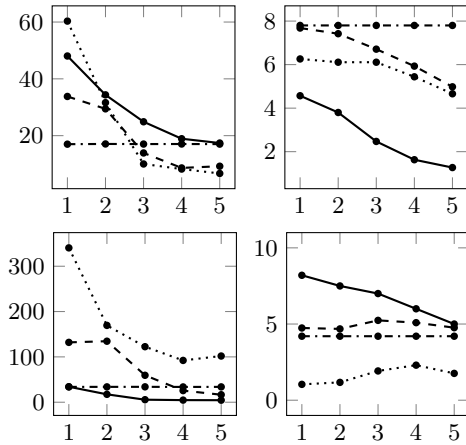


Figure 5: Local action-set property explanation at  $x = 1.0$  (top) and  $x = 1.5$  (bottom) SysS with trap learning. Average runtime (left) and average #MUGS (right) as a function of  $|A|$ . Blocksworld solid, NoMystery dashed, Rovers dash-dotted, TPP dotted.

Figure 5 shows performance data for local explanation. As before, this form of analysis becomes much easier as question size grows. runtime decreases significantly. Coverage (now shown in the figure to avoid cluttering the plots) increases accordingly. Explanation size, i. e., the number of MUGS, decreases as well (with the exception of Rovers as already pointed out).

## 7 Other Related Works

Various works relate to our computation of MUGS. In SAT and CP, minimal unsatisfiable cores (MUC, e. g. (Chinneck 2007; Laborie 2014)) are related to MUGS, yet rely on constraint-conflict analysis and are non-exhaustive, i. e., do not aim at finding *all* MUCs. Similarly, in the temporal planner IxTeT, an analysis of resource conflicts finds minimal critical sets (effects over-consuming a resource) (Laborie and Ghallab 1995). On the exhaustive side, problems and algorithms involving some form of solvability borderline within a lattice of problem variants have been considered since a long time (e. g. (de Kleer 1986; Reiter 1987)). In particular, the diagnosis framework by Grastien et al. (2011; 2012) can cast maximum solvable goal subsets as preferred diagnoses. That work uses SAT instead of a base planner, but an encoding into planning has been devised and ideas for search enhancements may be transferrable.

Some prior work (Yu et al. 2017; Lauffer and Topcu 2019) aims at suggesting goals to drop in oversubscribed situations, through conflict analyses related to finding a MUGS. At the technical level, both works are quite different from ours though as they use constraint-based problem encodings.

Prior work on other forms of plan-space explanation addressed unsolvable tasks, identifying minimal model changes resulting in solvability (Göbelbecker et al. 2010), or minimal differences between a solvable user model vs. unsolvable system model (Sreedharan et al. 2019).

We remark that our approach can be viewed as a form of domain/task analysis, or as a form of model checking applied to planning models. Both have been explored before, but addressing very different problems (e. g. (Fox and Long 1998; Rintanen 2000; Vaquero et al. 2013)).

Finally, to the reader versed in classical planning techniques, action-set properties may be reminiscent of disjunctive action landmarks (Karpas and Domshlak 2009; Helmert and Domshlak 2009), a technique used for computing heuristic functions. Disjunctive action landmarks correspond to atomic action-set properties. Propositional combinations other than conjunction (landmark collections) have not been considered though, and the previous purpose of disjunctive action landmarks is entirely different from ours.

## 8 Conclusion

Our approach analyses plan space in terms of plan properties and their dependencies. This naturally addresses challenges raised in previous work on iterative planning, and performance is reasonable in IPC benchmark studies. The approach can be easily generalized to almost arbitrary planning frameworks and plan-property languages, and we hope it will inspire other researchers too. To name just one example, plan properties and their dependencies might help in plan verbalization (Rosenthal, Selvaraj, and Veloso 2016), as an abstraction level (plan properties: “I did X”) as well as a means for justification (dependencies: “I did not do Y because ...”).

An interesting challenge is to answer deeper “why” questions regarding an entailment  $\Pi \models p \Rightarrow q$ . Possible ideas are to include additional plan properties elucidating the causal chain between  $p$  and  $q$ ; or to find a minimal relaxation (superset) of the plan set  $\Pi$  for which  $p$  no longer entails  $q$ , thus elucidating the circumstances under which that entailment holds. A computational idea worth trying is to represent the plan set  $\Pi$  symbolically and use that representation to identify entailment relations. We aim to apply and evaluate our approach in concrete use cases such as simulated penetration testing and AUV mission control (Cashmore et al. 2014).

## Acknowledgments

This material is based upon work supported by the Air Force Office of Scientific Research under award number FA9550-18-1-0245. Rebecca Eifler was also supported by the German Research Foundation (DFG) as part of CRC 248 (see [perspicuous-computing.science](http://perspicuous-computing.science)). Daniele Magazzeni’s group has also received support by EPSRC grant EP/R033722/1: Trust in Human-Machine Partnerships. Part of this work was performed while Jörg Hoffmann was visiting NASA Ames Research Center. We thank J. Benton, Minh Do, Jeremy Frank, and David Smith for insightful discussions.

## References

Bäckström, C., and Nebel, B. 1995. Complexity results for SAS<sup>+</sup> planning. *Computational Intelligence* 11(4):625–655.



- Baier, J. A.; Bacchus, F.; and McIlraith, S. A. 2009. A heuristic search approach to planning with temporally extended preferences. *AI* 173(5-6):593–618.
- Boddy, M.; Gohde, J.; Haigh, T.; and Harp, S. 2005. Course of action generation for cyber security using classical planning. In *Proc. ICAPS*, 12–21.
- Cashmore, M.; Fox, M.; Larkworthy, T.; Long, D.; and Magazzeni, D. 2014. AUV mission control via temporal planning. In *Proc. ICRA*, 6535–6541.
- Chinneck, J. W. 2007. *Feasibility and Infeasibility in Optimization: Algorithms and Computational Methods*. Springer-Verlag.
- de Kleer, J. 1986. An assumption-based TMS. *AI* 28(2):127–162.
- Domshlak, C., and Mirkis, V. 2015. Deterministic over-subscription planning as heuristic search: Abstractions and reformulations. *JAIR* 52:97–169.
- Edelkamp, S. 2006. On the compilation of plan constraints and preferences. In *Proc. ICAPS*, 374–377.
- Fox, M., and Long, D. 1998. The automatic inference of state invariants in TIM. *JAIR* 9:367–421.
- Fox, M.; Long, D.; and Magazzeni, D. 2017. Explainable planning. In *Proc. IJCAI’17 Workshop on Explainable AI*.
- Gazen, B. C., and Knoblock, C. 1997. Combining the expressiveness of UCPOP with the efficiency of Graphplan. In *Proc. ECP*, 221–233.
- Göbelbecker, M.; Keller, T.; Eyerich, P.; Brenner, M.; and Nebel, B. 2010. Coming up with good excuses: What to do when no plan can be found. In *Proc. ICAPS*, 81–88.
- Grastien, A.; Haslum, P.; and Thiébaux, S. 2011. Exhaustive diagnosis of discrete event systems through exploration of the hypothesis space. In *Proc. International Workshop on Principles of Diagnosis (DX)*.
- Grastien, A.; Haslum, P.; and Thiébaux, S. 2012. Conflict-based diagnosis of discrete event systems: Theory and practice. AAAI Press.
- Helmert, M., and Domshlak, C. 2009. Landmarks, critical paths and abstractions: What’s the difference anyway? In *Proc. ICAPS*, 162–169.
- Helmert, M. 2006. The Fast Downward planning system. *JAIR* 26:191–246.
- Helmert, M. 2009. Concise finite-domain representations for PDDL planning tasks. *AI* 173:503–535.
- Hoffmann, J., and Nebel, B. 2001. The FF planning system: Fast plan generation through heuristic search. *JAIR* 14:253–302.
- Hoffmann, J. 2015. Simulated penetration testing: From “Dijkstra” to “Turing Test++”. In *Proc. ICAPS*.
- Karpas, E., and Domshlak, C. 2009. Cost-optimal planning with landmarks. In Boutilier, C., ed., *Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI’09)*, 1728–1733. Pasadena, California, USA: Morgan Kaufmann.
- Katz, M.; Keyder, E.; Winterer, D.; and Pommerening, F. 2019. Oversubscription planning as classical planning with multiple cost functions. In *Proc. ICAPS*, 237–245.
- Laborie, P., and Ghallab, M. 1995. Planning with sharable resource constraints. In *Proc. IJCAI*, 1643–1649.
- Laborie, P. 2014. An optimal iterative algorithm for extracting mucs in a black-box constraint network. In *Proc. ECAI*, 1051–1052.
- Lauffer, N., and Topcu, U. 2019. Human-understandable explanations of infeasibility for resource-constrained scheduling problems. In *Proc. 2nd Workshop on Explainable Planning (XAIP’19)*.
- Miller, T. 2019. Explanation in artificial intelligence: Insights from the social sciences. *AI* 267:1–38.
- Muise, C. 2016. Planning domains. In *ICAPS System Demonstrations and Exhibits*.
- Nakhost, H.; Hoffmann, J.; and Müller, M. 2012. Resource-constrained planning: A Monte Carlo random walk approach. In *Proc. ICAPS*, 181–189.
- Nebel, B. 2000. On the compilability and expressive power of propositional planning formalisms. *JAIR* 12:271–315.
- Palacios, H., and Geffner, H. 2009. Compiling uncertainty away in conformant planning problems with bounded width. *JAIR* 35:623–675.
- Reiter, R. 1987. A theory of diagnosis from first principles. *AI* 32(1):57–95.
- Rintanen, J. 2000. An iterative algorithm for synthesizing invariants. In *Proc. AAAI*, 806–811.
- Rosenthal, S.; Selvaraj, S. P.; and Veloso, M. M. 2016. Verbalization: Narration of autonomous robot experience. In *Proc. IJCAI*.
- Smith, D. E. 2004. Choosing objectives in over-subscription planning. In *Proc. ICAPS*, 393–401.
- Smith, D. 2012. Planning as an iterative process. In *Proc. AAAI*, 2180–2185.
- Sreedharan, S.; Srivastava, S.; Smith, D.; and Kambhampati, S. 2019. Why couldn’t you do that? explaining unsolvability of classical planning problems in the presence of plan advice. In *Proc. IJCAI*.
- Steinmetz, M., and Hoffmann, J. 2017a. Search and learn: On dead-end detectors, the traps they set, and trap learning. In *Proc. IJCAI*.
- Steinmetz, M., and Hoffmann, J. 2017b. State space search nogood learning: Online refinement of critical-path dead-end detectors in planning. *AI* 245:1–37.
- Vaquero, T. S.; Silva, J. R.; Tonidandel, F.; and Beck, J. C. 2013. itsimple: towards an integrated design system for real planning applications. *Knowledge Engineering Review* 28(2):215–230.
- Yu, P.; Williams, B. C.; Fang, C.; Cui, J.; and Haslum, P. 2017. Resolving over-constrained temporal problems with uncertainty through conflict-directed relaxation. *JAIR* 60:425–490.