

# Learning STRIPS action models with classical planning

Diego Aineto and Sergio Jiménez and Eva Onaindia

Departamento de Sistemas Informáticos y Computación

Universitat Politècnica de València.

Camino de Vera s/n. 46022 Valencia, Spain

{dieaigar,serjice,onaindia}@dsic.upv.es

## Abstract

This paper presents a novel approach for learning STRIPS action models from examples that compiles this inductive learning task into classical planning. Interestingly, the compilation approach is flexible to different amounts of available input knowledge; the learning examples can range from a set of plans (with their corresponding initial and final states) to just a set of initial and final states (no intermediate action or state is given). What is more, the compilation accepts partially specified action models and can be used to validate whether the observation of a plan execution follows a given STRIPS action model, even if this model is not fully specified.

## Introduction

Besides *plan synthesis* (Ghallab, Nau, and Traverso 2004), planning action models are also useful for *plan/goal recognition* (Ramírez 2012). At both planning tasks, an automated planner is required to reason about action models that correctly and completely capture the possible world transitions (Geffner and Bonet 2013). Unfortunately, building planning action models is complex, even for planning experts, and this knowledge acquisition task is a bottleneck that limits the potential of AI planning (Kambhampati 2007).

On the other hand, Machine Learning (ML) has shown to be able to compute a wide range of different kinds of models from examples (Michalski, Carbonell, and Mitchell 2013). The application of inductive ML to the learning of STRIPS action models, the vanilla action model for planning (Fikes and Nilsson 1971), is not straightforward though:

- The *input* to ML algorithms (the learning/training data) usually are finite vectors encoding the value of fixed features in a given set of objects. The input for learning planning action models are observations of plan executions (where each plan possibly has a different length).
- The *output* of ML algorithms usually is a scalar value (an integer, in the case of *classification* tasks, or a real value, in the case of *regression* tasks). When learning STRIPS action models the output is, for each action, the sets of preconditions, negative and positive effects, that define the possible state transitions.

Learning STRIPS action models is a well-studied problem with sophisticated algorithms, like ARMS (Yang, Wu, and Jiang 2007), SLAF (Amir and Chang 2008) or LOCM (Cresswell, McCluskey, and West 2013) that do not require full knowledge of the intermediate states traversed by the example plans. Motivated by recent advances on the synthesis of different kinds of generative models with classical planning (Bonet, Palacios, and Geffner 2009; Segovia-Aguas, Jiménez, and Jonsson 2016; 2017), this paper introduces an innovative approach for learning STRIPS action models that can be defined as a classical planning compilation. The compilation approach is appealing by itself because opens the door to the bootstrapping of planning action models but also because:

1. Is flexible to different amounts of available input knowledge. The learning examples can range from a set of plans (with their corresponding initial and final states) to just a set of initial and final states where no intermediate state or action is observed.
2. Accepts previous knowledge about the structure of the actions in the form of partially specified action models. In the extreme, the compilation can validate whether an observed plan execution is valid for a given STRIPS action model, even if this model is not fully specified.

The second section of the paper formalizes the classical planning model, its extension to *conditional effects* (a requirement of the proposed compilation) and the STRIPS action model (the output of the addressed learning task). The third section formalizes the learning of STRIPS action models with regard to different amounts of available input knowledge. The fourth and fifth sections describe our approach for addressing the formalized learning tasks. Finally, the last sections report the data collected in a empirical evaluation, discuss the strengths and weaknesses of our approach and propose several opportunities for future research.

## Background

This section defines the planning models used on this work and the output of the learning tasks addressed in the paper.

### Classical planning

We use  $F$  to denote the set of *fluents* (propositional variables) describing a state. A *literal*  $l$  is a valuation of a fluent

$f \in F$ , i.e. either  $l = f$  or  $l = \neg f$ . A set of literals  $L$  represents a partial assignment of values to fluents (WLOG we assume that  $L$  does not assign conflicting values to any fluent). We use  $\mathcal{L}(F)$  to denote the set of all literal sets on  $F$ , i.e. all partial assignments of values to fluents.

A *state*  $s$  is a full assignment of values to fluents, i.e.  $|s| = |F|$ , so the size of the state space is  $2^{|F|}$ . Explicitly including negative literals  $\neg f$  in states simplifies subsequent definitions but often, we will abuse notation by defining a state  $s$  only in terms of the fluents that are true in  $s$ , as is common in STRIPS planning.

A *classical planning frame* is a tuple  $\Phi = \langle F, A \rangle$ , where  $F$  is a set of fluents and  $A$  is a set of actions. An action  $a \in A$  is defined with *preconditions*,  $\text{pre}(a) \subseteq \mathcal{L}(F)$ , *positive effects*,  $\text{eff}^+(a) \subseteq \mathcal{L}(F)$ , and *negative effects*  $\text{eff}^-(a) \subseteq \mathcal{L}(F)$ . We say that an action  $a \in A$  is *applicable* in a state  $s$  iff  $\text{pre}(a) \subseteq s$ . The result of applying  $a$  in  $s$  is the *successor state* denoted by  $\theta(s, a) = \{s \setminus \text{eff}^-(a)\} \cup \text{eff}^+(a)$ .

We say that an action  $a \in A$  is *applicable* in a state  $s$  iff  $\text{pre}(a) \subseteq s$ . The result of applying  $a$  in  $s$  is the *successor state*  $\theta(s, a) = \{s \setminus \text{eff}^-(a)\} \cup \text{eff}^+(a)$ .

A *classical planning problem* is a tuple  $P = \langle F, A, I, G \rangle$ , where  $I$  is an initial state and  $G \subseteq \mathcal{L}(F)$  is a goal condition. A *plan* for  $P$  is an action sequence  $\pi = \langle a_1, \dots, a_n \rangle$  that induces a state sequence  $\langle s_0, s_1, \dots, s_n \rangle$  such that  $s_0 = I$  and, for each  $1 \leq i \leq n$ ,  $a_i$  is applicable in  $s_{i-1}$  and generates the successor state  $s_i = \theta(s_{i-1}, a_i)$ . We denote with  $|\pi|$  the *plan length*. A plan  $\pi$  *solves*  $P$  iff  $G \subseteq s_n$ , i.e. if the goal condition is satisfied at the last state reached after following the application of  $\pi$  in  $I$ .

### Classical planning with conditional effects

Our approach for learning STRIPS action models is compiling this learning task into a classical planning task with conditional effects. Conditional effects allow us to compactly define actions whose effects depend on the current state. Supporting conditional effects is now a requirement of the IPC (Vallati et al. 2015) and many classical planners cope with conditional effects without compiling them away.

An action  $a \in A$  has now a set of *preconditions*  $\text{pre}(a) \in \mathcal{L}(F)$  and a set of *conditional effects*  $\text{cond}(a)$ . Each conditional effect  $C \triangleright E \in \text{cond}(a)$  is composed of two sets of literals  $C \in \mathcal{L}(F)$ , the *condition*, and  $E \in \mathcal{L}(F)$ , the *effect*.

An action  $a \in A$  is *applicable* in a state  $s$  if and only if  $\text{pre}(a) \subseteq s$ , and the resulting set of *triggered effects* are the effects whose conditions hold in  $s$ :

$$\text{triggered}(s, a) = \bigcup_{C \triangleright E \in \text{cond}(a), C \subseteq s} E,$$

The result of applying an action  $a$  in a state  $s$  is the *successor state*  $\theta(s, a) = \{s \setminus \text{eff}_c^-(s, a)\} \cup \text{eff}_c^+(s, a)$  where  $\text{eff}_c^-(s, a) \subseteq \text{triggered}(s, a)$  and  $\text{eff}_c^+(s, a) \subseteq \text{triggered}(s, a)$  are the triggered *negative* and *positive* effects, respectively.

### STRIPS action schemes and variable name objects

This work addresses the learning of PDDL action schemes that follow the STRIPS requirement (McDermott et al.

```
(:action stack
:parameters (?v1 ?v2 - object)
:precondition (and (holding ?v1) (clear ?v2))
:effect (and (not (holding ?v1))
(not (clear ?v2))
(handempty) (clear ?v1)
(on ?v1 ?v2)))
```

Figure 1: STRIPS operator schema coding, in PDDL, the *stack* action from the *blocksworld*.

1998; Fox and Long 2003). Figure 1 shows the schema, coded in PDDL, for the *stack* action from a four-operator *blocksworld* (Slaney and Thiébaux 2001).

To formalize the output of the learning task, we assume that fluents  $F$  are instantiated from a set of *predicates*  $\Psi$ , as in PDDL. Each predicate  $p \in \Psi$  has an argument list of arity  $\text{ar}(p)$ . Given a set of *objects*  $\Omega$ , the set of fluents  $F$  is induced by assigning objects in  $\Omega$  to the arguments of predicates in  $\Psi$ , i.e.  $F = \{p(\omega) : p \in \Psi, \omega \in \Omega^{\text{ar}(p)}\}$  s.t.  $\Omega^k$  is the  $k$ -th Cartesian power of  $\Omega$ .

Let  $\Omega_v = \{v_i\}_{i=1}^{\max_{a \in A} \text{ar}(a)}$  be a new set of objects  $\Omega \cap \Omega_v = \emptyset$ , denoted as *variable names*, and that is bound by the maximum arity of an action in a given planning frame. For instance, in a three-block *blocksworld*  $\Omega = \{\text{block}_1, \text{block}_2, \text{block}_3\}$  while  $\Omega_v = \{v_1, v_2\}$  because the operators with the maximum arity, *stack* and *unstack*, have two parameters each.

Let us also define  $F_v$ , a new set of fluents  $F \cap F_v = \emptyset$ , that results from instantiating  $\Psi$  using only the objects in  $\Omega_v$  and that defines the elements that can appear in an action schema. For instance, in the *blocksworld*,  $F_v = \{\text{handempty}, \text{holding}(v_1), \text{holding}(v_2), \text{clear}(v_1), \text{clear}(v_2), \text{ontable}(v_1), \text{ontable}(v_2), \text{on}(v_1, v_1), \text{on}(v_1, v_2), \text{on}(v_2, v_1), \text{on}(v_2, v_2)\}$ .

Finally, we assume that actions  $a \in A$  are instantiated from STRIPS operator schemes  $\xi = \langle \text{head}(\xi), \text{pre}(\xi), \text{add}(\xi), \text{del}(\xi) \rangle$  where:

- $\text{head}(\xi) = \langle \text{name}(\xi), \text{pars}(\xi) \rangle$ , is the operator *header* defined by its name and corresponding *variable names*,  $\text{pars}(\xi) = \{v_i\}_{i=1}^{\text{ar}(\xi)}$ . For instance, the headers for a four-operator *blocksworld* are:  $\text{pickup}(v_1)$ ,  $\text{putdown}(v_1)$ ,  $\text{stack}(v_1, v_2)$  and  $\text{unstack}(v_1, v_2)$ .
- The preconditions  $\text{pre}(\xi) \subseteq F_v$ , the negative effects  $\text{del}(\xi) \subseteq F_v$ , and the positive effects  $\text{add}(\xi) \subseteq F_v$  such that,  $\text{del}(\xi) \subseteq \text{pre}(\xi)$ ,  $\text{del}(\xi) \cap \text{add}(\xi) = \emptyset$  and  $\text{pre}(\xi) \cap \text{add}(\xi) = \emptyset$ .

### Learning STRIPS action models

Learning STRIPS action models from fully available input knowledge, i.e. from plans where the *pre*- and *post*-states of every action in a plan are available, is straightforward. When any intermediate state is available, STRIPS operator schemes are derived lifting the literals that change between the pre and post-state of the corresponding action executions. Preconditions are derived lifting the minimal set of literals that

appears in all the pre-states of the corresponding actions.

This section formalizes more challenging learning tasks, where less input knowledge is available:

**Learning from (initial, final) state pairs.** This learning task corresponds to observing an agent acting in the world but watching only the results of its plan executions. No information about the actions in the plans is given. This learning task is formalized as  $\Lambda = \langle \Psi, \Sigma \rangle$ :

- $\Psi$  is the set of predicates that define the abstract state space of a given planning domain.
- $\Sigma = \{\sigma_1, \dots, \sigma_\tau\}$  is a set of (initial, final) state pairs, that we call *labels*. Each label  $\sigma_t = (s_0^t, s_n^t)$ ,  $1 \leq t \leq \tau$ , comprises the *final* state  $s_n^t$  resulting from executing an unknown plan  $\pi_t$  starting from the *initial* state  $s_0^t$ .

**Learning from labeled plans.** Here we augment the input knowledge with the actions executed by the observed agent and define the learning task  $\Lambda' = \langle \Psi, \Sigma, \Pi \rangle$ :

- $\Pi = \{\pi_1, \dots, \pi_\tau\}$  is a given set of example plans where each plan  $\pi_t = \langle a_1^t, \dots, a_n^t \rangle$ ,  $1 \leq t \leq \tau$ , is an action sequence that induces the corresponding state sequence  $\langle s_0^t, s_1^t, \dots, s_n^t \rangle$  such that, for each  $1 \leq i \leq n$ ,  $a_i^t$  is applicable in  $s_{i-1}^t$  and generates  $s_i^t = \theta(s_{i-1}^t, a_i^t)$ .

Figure 2 shows an example of a learning task  $\Lambda'$  in the blockworld. This learning task has a single learning example,  $\Pi = \{\pi_1\}$  and  $\Sigma = \{\sigma_1\}$ , that corresponds to observing the execution of an eight-action plan ( $|\pi_1| = 8$ ) for inverting a four-block tower.

**Learning from partially specified action models.** We may not require to start learning from scratch so we augment the learning input with partially specified operator schemes. This learning task is defined as  $\Lambda'' = \langle \Psi, \Sigma, \Pi, \Xi_0 \rangle$ :

- $\Xi_0$  is a partially specified action model in which some preconditions and effects are a priori known.

A solution to  $\Lambda$  is a set of operator schema  $\Xi$  that is compliant just with the predicates in  $\Psi$ , and the given set of initial and final states  $\Sigma$ . In this learning scenario, a solution must not only determine a possible STRIPS action model but also the plans  $\pi_t$ ,  $1 \leq t \leq \tau$  that explain the given labels  $\Sigma$  using the learned STRIPS model. A solution to  $\Lambda'$  is a set of STRIPS operator schema  $\Xi$  (one schema  $\xi = \langle head(\xi), pre(\xi), add(\xi), del(\xi) \rangle$  for each action with a different name in the example plans  $\Pi$ ) compliant with the predicates in  $\Psi$ , the example plans  $\Pi$ , and their corresponding labels  $\Sigma$ . Finally a solution to  $\Lambda''$  is a set of STRIPS operator schema  $\Xi$  compliant as well with the provided partially specified action model  $\Xi_0$ .

## Learning STRIPS action models with planning

Our approach for addressing a learning task  $\Lambda$ ,  $\Lambda'$  or  $\Lambda''$ , is compiling it into a classical planning task with conditional effects. The intuition behind the compilation is that a solution to the resulting classical planning task is a sequence of actions that:

;;; Predicates in  $\Psi$

```
(handempty) (holding ?o - object)
(clear ?o - object) (ontable ?o - object)
(on ?o1 - object ?o2 - object)
```

;;; Plan  $\pi_1$

```
0: (unstack A B)
1: (putdown A)
2: (unstack B C)
3: (stack B A)
4: (unstack C D)
5: (stack C B)
6: (pickup D)
7: (stack D C)
```

;;; Label  $\sigma_1 = (s_0^1, s_n^1)$

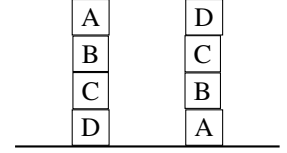


Figure 2: Example of a task for learning a STRIPS action model in the blockworld from a single labeled plan.

1. Programs the STRIPS action model  $\Xi$ . A solution plan has a *prefix* that, for each  $\xi \in \Xi$ , determines the fluents from  $F_v$  that belong to its *pre*( $\xi$ ), *del*( $\xi$ ) and *add*( $\xi$ ) sets.
2. Validates the programmed STRIPS action model  $\Xi$  in the given input knowledge (the labels  $\Sigma$ , and  $\Pi$  and/or  $\Xi_0$  if available). For every label  $\sigma_t \in \Sigma$ , a solution plan has a postfix that produces a final state  $s_n^t$  starting from the corresponding initial state  $s_0^t$  using the programmed action model  $\Xi$ . We call this process the validation of the programmed STRIPS action model  $\Xi$ , at the learning example  $1 \leq t \leq \tau$ .

To formalize our compilation we first define  $1 \leq t \leq \tau$  classical planning instances  $P_t = \langle F, \emptyset, I_t, G_t \rangle$  that belong to the same planning frame (i.e. same fluents and actions but differ in the initial state and goals). Fluents  $F$  are built instantiating the predicates in  $\Psi$  with the objects appearing in the input labels  $\Sigma$ . Formally  $\Omega = \{o|o \in \bigcup_{1 \leq t \leq \tau} obj(s_0^t)\}$ , where *obj* is a function that returns the set of objects that appear in a fully specified state. The set of actions,  $A = \emptyset$ , is empty because the action model is initially unknown. Finally, the initial state  $I_t$  is given by the state  $s_0^t \in \sigma_t$  while goals  $G_t$ , are defined by the state  $s_n^t \in \sigma_t$ .

Now we are ready to formalize the compilations. We start with  $\Lambda$ , because it requires less input knowledge. Given a learning task  $\Lambda = \langle \Psi, \Sigma \rangle$  the compilation outputs a classical planning task  $P_\Lambda = \langle F_\Lambda, A_\Lambda, I_\Lambda, G_\Lambda \rangle$ :

- $F_\Lambda$  extends  $F$  with:
  - Fluents representing the programmed action model  $pre_f(\xi)$ ,  $del_f(\xi)$  and  $add_f(\xi)$ , for every  $f \in F_v$  and  $\xi \in \Xi$ . If a fluent  $pre_f(\xi)/del_f(\xi)/add_f(\xi)$  holds, it means that  $f$  is a precondition/negative effect/positive effect in the STRIPS operator schema  $\xi \in \Xi$ . For instance, the preconditions of the *stack* schema (Figure 1) are represented by fluents *pre\_holding\_stack.v1* and *pre\_clear\_stack.v2*.
  - A fluent *mode<sub>prog</sub>* indicating whether the operator schemes are being programmed or validated (already programmed) and fluents  $\{test_t\}_{1 \leq t \leq \tau}$ , indicating the example where the action model is being validated.

- $I_\Lambda$  contains the fluents from  $F$  that encode  $s_0^1$  (the initial state of the first label), every  $pre_f(\xi) \in F_\Lambda$  and  $mode_{prog}$  set to true. Our compilation assumes that initially any operator schema is programmed with every possible precondition, no negative effect and no positive effect.
- $G_\Lambda = \bigcup_{1 \leq t \leq \tau} \{test_t\}$ , indicates that the programmed action model is validated in all the learning examples.
- $A_\Lambda$  contains actions of three kinds:
  1. Actions for *programming* an operator schema  $\xi \in \Xi$ :
    - Actions for **removing** a *precondition*  $f \in F_v$  from the action schema  $\xi \in \Xi$ .

$$\begin{aligned} pre(\text{programPre}_{f,\xi}) &= \{\neg del_f(\xi), \neg add_f(\xi), \\ &\quad mode_{prog}, pre_f(\xi)\}, \\ cond(\text{programPre}_{f,\xi}) &= \{\emptyset\} \triangleright \{\neg pre_f(\xi)\}. \end{aligned}$$

- Actions for **adding** a *negative* or *positive* effect  $f \in F_v$  to the action schema  $\xi \in \Xi$ .

$$\begin{aligned} pre(\text{programEff}_{f,\xi}) &= \{\neg del_f(\xi), \neg add_f(\xi), \\ &\quad mode_{prog}\}, \\ cond(\text{programEff}_{f,\xi}) &= \{pre_f(\xi)\} \triangleright \{del_f(\xi)\}, \\ &\quad \{\neg pre_f(\xi)\} \triangleright \{add_f(\xi)\}. \end{aligned}$$

2. Actions for *applying* an already programmed operator schema  $\xi \in \Xi$  bound with the objects  $\omega \subseteq \Omega^{ar(\xi)}$ . We assume operators headers are known so the binding of the operator schema is done implicitly by order of appearance of the action parameters, i.e. variables  $pars(\xi)$  are bound to the objects in  $\omega$  appearing at the same position. Figure 3 shows the PDDL encoding of the action for applying a programmed operator *stack*.

$$\begin{aligned} pre(\text{apply}_{\xi,\omega}) &= \{pre_f(\xi) \implies p(\omega)\}_{\forall p \in \Psi, f=p(pars(\xi))}, \\ cond(\text{apply}_{\xi,\omega}) &= \{del_f(\xi)\} \triangleright \{\neg p(\omega)\}_{\forall p \in \Psi, f=p(pars(\xi))}, \\ &\quad \{add_f(\xi)\} \triangleright \{p(\omega)\}_{\forall p \in \Psi, f=p(pars(\xi))}, \\ &\quad \{mode_{prog}\} \triangleright \{\neg mode_{prog}\}. \end{aligned}$$

3. Actions for *validating* the learning example  $1 \leq t \leq \tau$ .
 
$$\begin{aligned} pre(\text{validate}_t) &= G_t \cup \{test_j\}_{j \in 1 \leq j < t} \\ &\quad \cup \{\neg test_j\}_{j \in t \leq j \leq \tau} \cup \{\neg mode_{prog}\}, \\ cond(\text{validate}_t) &= \{\emptyset\} \triangleright \{test_t\}. \end{aligned}$$

**Lemma 1.** Any classical plan  $\pi$  that solves  $P_\Lambda$  induces an action model  $\Xi$  that solves the learning task  $\Lambda$ .

*Proof sketch.* The compilation forces that once the preconditions of an operator schema  $\xi \in \Xi$  are programmed, they cannot be altered. The same happens with the positive and negative effects that define an operator schema  $\xi \in \Xi$  (besides they can only be programmed after preconditions are programmed). Once operator schemas are programmed they can only be applied because of the  $mode_{prog}$  fluent. To solve  $P_\Lambda$ , goals  $\{test_t\}$ ,  $1 \leq t \leq \tau$  can only be achieved: executing an applicable sequence of programmed operator schemes that reaches the final state  $s_n^t$ , defined in  $\sigma_t$ , starting from  $s_0^t$ . If this is achieved for all the input examples  $1 \leq t \leq \tau$ , it means that the programmed action model  $\Xi$  is compliant with the provided input knowledge and hence, it is a solution to  $\Lambda$ .  $\square$

```
(:action apply_stack
:parameters (?o1 - object ?o2 - object)
:precondition
  (and (or (not (pre_on_stack_v1_v1)) (on ?o1 ?o1))
        (or (not (pre_on_stack_v1_v2)) (on ?o1 ?o2))
        (or (not (pre_on_stack_v2_v1)) (on ?o2 ?o1))
        (or (not (pre_on_stack_v2_v2)) (on ?o2 ?o2))
        (or (not (pre_ontable_stack_v1)) (ontable ?o1))
        (or (not (pre_ontable_stack_v2)) (ontable ?o2))
        (or (not (pre_clear_stack_v1)) (clear ?o1))
        (or (not (pre_clear_stack_v2)) (clear ?o2))
        (or (not (pre_holding_stack_v1)) (holding ?o1))
        (or (not (pre_holding_stack_v2)) (holding ?o2))
        (or (not (pre_handempty_stack)) (handempty)))
:effect
  (and (when (del_on_stack_v1_v1) (not (on ?o1 ?o1)))
        (when (del_on_stack_v1_v2) (not (on ?o1 ?o2)))
        (when (del_on_stack_v2_v1) (not (on ?o2 ?o1)))
        (when (del_on_stack_v2_v2) (not (on ?o2 ?o2)))
        (when (del_ontable_stack_v1) (not (ontable ?o1)))
        (when (del_ontable_stack_v2) (not (ontable ?o2)))
        (when (del_clear_stack_v1) (not (clear ?o1)))
        (when (del_clear_stack_v2) (not (clear ?o2)))
        (when (del_holding_stack_v1) (not (holding ?o1)))
        (when (del_holding_stack_v2) (not (holding ?o2)))
        (when (del_handempty_stack) (not (handempty)))
        (when (add_on_stack_v1_v1) (on ?o1 ?o1))
        (when (add_on_stack_v1_v2) (on ?o1 ?o2))
        (when (add_on_stack_v2_v1) (on ?o2 ?o1))
        (when (add_on_stack_v2_v2) (on ?o2 ?o2))
        (when (add_ontable_stack_v1) (ontable ?o1))
        (when (add_ontable_stack_v2) (ontable ?o2))
        (when (add_clear_stack_v1) (clear ?o1))
        (when (add_clear_stack_v2) (clear ?o2))
        (when (add_holding_stack_v1) (holding ?o1))
        (when (add_holding_stack_v2) (holding ?o2))
        (when (add_handempty_stack) (handempty))
        (when (modeProg) (not (modeProg)))))
```

Figure 3: Action for applying an already programmed schema *stack* as encoded in PDDL (implications coded as disjunctions).

The compilation is *complete* in the sense that it does not discard any possible STRIPS action model.

## Constraining the learning hypothesis space with additional input knowledge

Here we show that further input knowledge can be used to constrain the space of possible action models and make the learning of STRIPS action models more practicable.

### Labeled plans

We extend the compilation to consider labeled plans. Given a learning task  $\Lambda' = \langle \Psi, \Sigma, \Pi \rangle$ , the compilation outputs a classical planning task  $P_{\Lambda'} = \langle F_{\Lambda'}, A_{\Lambda'}, I_{\Lambda'}, G_{\Lambda'} \rangle$  that extends  $P_\Lambda$  as follows:

- $F_{\Lambda'}$  extends  $F_\Lambda$  with  $F_\Pi = \{plan(name(\xi), \Omega^{ar(\xi)}, j)\}$ , the fluents to code the steps of the plans in  $\Pi$ , where  $F_{\pi_t} \subseteq F_\Pi$  encodes  $\pi_t \in \Pi$ . Fluents  $at_j$  and  $next_{j,j_2}$ ,  $1 \leq j < j_2 \leq n$ , are also added to represent the current plan step and to iterate through the steps of a plan.

- $I_{\Lambda'}$  extends  $I_{\Lambda}$  with fluents  $F_{\pi_1}$  plus fluents  $at_1$  and  $\{next_{j,j_2}\}$ ,  $1 \leq j < j_2 \leq n$ , for indicating the plan step where the action model is validated. Goals are  $G_{\Lambda'} = G_{\Lambda} = \bigcup_{1 \leq t \leq \tau} \{test_t\}$ , as in the original compilation.
- With respect to  $A_{\Lambda'}$ .
  1. The actions for *programming* the preconditions/effects of a given operator schema  $\xi \in \Xi$  are the same.
  2. The actions for *applying* an already programmed operator have an extra precondition  $f \in F_{\Pi}$ , that encodes the current plan step, and extra conditional effects  $\{at_j\} \triangleright \{\neg at_j, at_{j+1}\}_{\forall j \in [1,n]}$  for advancing to the next plan step. This mechanism forces that these actions are only applied as in the example plans.
  3. The actions for *validating* the current learning example have an extra precondition,  $at_{|\pi_t|}$ , to indicate that the current plan  $\pi_t$  was fully executed and extra conditional effects to unload plan  $\pi_t$  and load the next plan  $\pi_{t+1}$ :
$$\{\emptyset\} \triangleright \{\neg at_{|\pi_t|}, at_1\}, \{f\} \triangleright \{\neg f\}_{f \in F_{\pi_t}}, \{\emptyset\} \triangleright \{f\}_{f \in F_{\pi_{t+1}}}.$$

### Partially specified action models

Known preconditions and effects are encoded as fluents  $pre_f(\xi)$ ,  $del_f(\xi)$  and  $add_f(\xi)$  set to true at the initial state  $I_{\Lambda'}$ . The corresponding programming actions,  $programPref_{f,\xi}$  and  $programEff_{f,\xi}$ , become unnecessary and are removed from  $A_{\Lambda'}$  making the classical planning task  $P_{\Lambda'}$  easier to be solved.

To illustrate this, the classical plan of Figure 4 is a solution to a learning task  $\Lambda'' = \langle \Psi, \Sigma, \Pi, \Xi_0 \rangle$  for getting the blocksworld action model where operator schemes for *pickup*, *putdown* and *unstack* are specified in  $\Xi_0$ . This plan programs and validates the operator schema *stack* from the blocksworld, using the plan  $\pi_1$  and label  $\sigma_1$  shown in Figure 2. Plan steps [0, 8] are the actions for programming the preconditions of the *stack* operator, steps [9, 13] are the actions for programming the operator effects and steps [14, 22] are the actions for validating the programmed operators following the plan  $\pi_1$  shown in the Figure 2.

In the extreme, when a fully specified STRIPS action model  $\Xi$  is given in  $\Xi_0$ , the compilation validates whether an observed plan follows the given model. In this case, if a solution plan is found to  $P_{\Lambda'}$ , it means that the given STRIPS action model is *valid* for the given examples. If  $P_{\Lambda'}$  is unsolvable it means that the given STRIPS action model is invalid since it is not compliant with all the given examples. Tools for plan validation like VAL (Howey, Long, and Fox 2004) could also be used at this point.

### Static predicates

A *static predicate*  $p \in \Psi$  is a predicate that does not appear in the effects of any action schema (Fox and Long 1998). Therefore, one can get rid of the mechanism for programming these predicates in the effects of any action schema while keeping the compilation complete. Formally, given a static predicate  $p$ :

- Fluents  $del_f(\xi)$  and  $add_f(\xi)$ , such that  $f \in F_v$  is an instantiation of the static predicate  $p$  in the set of *variable objects*  $\Omega_v$ , can be discarded for every  $\xi \in \Xi$ .

```

00 : (program_pre_clear_stack.v1)
01 : (program_pre_handempty_stack)
02 : (program_pre_holding_stack.v2)
03 : (program_pre_on_stack.v1.v1)
04 : (program_pre_on_stack.v1.v2)
05 : (program_pre_on_stack.v2.v1)
06 : (program_pre_on_stack.v2.v2)
07 : (program_pre_ontable_stack.v1)
08 : (program_pre_ontable_stack.v2)
09 : (program_eff_clear_stack.v1)
10 : (program_eff_clear_stack.v2)
11 : (program_eff_handempty_stack)
12 : (program_eff_holding_stack.v1)
13 : (program_eff_on_stack.v1.v2)
14 : (apply_unstack a b i1 i2)
15 : (apply_putdown a i2 i3)
16 : (apply_unstack b c i3 i4)
17 : (apply_stack b a i4 i5)
18 : (apply_unstack c d i5 i6)
19 : (apply_stack c b i6 i7)
20 : (apply_pickup d i7 i8)
21 : (apply_stack d c i8 i9)
22 : (validate_1)

```

Figure 4: Plan for programming and validating the *stack* schema using plan  $\pi_1$  and label  $\sigma_1$  (shown in Figure 2) as well as previously specified operator schemes for *pickup*, *putdown* and *unstack*.

- Actions  $programEff_{f,\xi}$  (s.t.  $f \in F_v$  is an instantiation of  $p$  in  $\Omega_v$ ) can also be discarded for every  $\xi \in \Xi$ .

Static predicates can also constrain the space of possible preconditions by looking at the given set of labels  $\Sigma$ . One can assume that if a precondition  $f \in F_v$  (s.t.  $f \in F_v$  is an instantiation of a static predicate in  $\Omega_v$ ) is not compliant with the labels in  $\Sigma$  then, fluents  $pre_f(\xi)$  and actions  $programPref_{f,\xi}$  can be discarded for every  $\xi \in \Xi$ . For instance in the *zenotravel* domain  $pre\_next\_board.v1.v1$ ,  $pre\_next\_debark.v1.v1$ ,  $pre\_next\_fly.v1.v1$ ,  $pre\_next\_zoom.v1.v1$ ,  $pre\_next\_refuel.v1.v1$  can be discarded (and their corresponding programming actions) because a precondition (next ?v1 ?v1 - flevel) will never hold at any state in  $\Sigma$ .

Likewise, static predicates can constrain the space of possible preconditions looking at the given example plans. Fluents  $pre_f(\xi)$  and actions  $programPref_{f,\xi}$  are discardable for every  $\xi \in \Xi$  if a precondition  $f \in F_v$  (s.t.  $f \in F_v$  is an instantiation of a static predicate in  $\Omega_v$ ) is not possible according to  $\Pi$ . Back to the *zenotravel* domain, if a example plan  $\pi_t \in \Pi$  contains the action (fly plane1 city2 city0 f13 f12) and the corresponding label  $\sigma_t \in \Sigma$  contains the static literal (next f12 f13) but does not contain (next f12 f12), (next f13 f13) or (next f13 f12) the only possible precondition including the static predicate is  $pre\_next\_fly.v5.v4$ .

In the evaluation of our approach we do not assume that the set of static predicates is given but compute a set of *potential* static predicates from the input to the learning task.

|            | No Static |      |      |      |      |      |      |      | Static |      |      |      |      |      |      |      |
|------------|-----------|------|------|------|------|------|------|------|--------|------|------|------|------|------|------|------|
|            | Pre       |      | Add  |      | Del  |      |      |      | Pre    |      | Add  |      | Del  |      |      |      |
|            | P         | R    | P    | R    | P    | R    | P    | R    | P      | R    | P    | R    | P    | R    | P    | R    |
| Blocks     | 1.0       | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0    | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  |
| Driverlog  | 1.0       | 0.36 | 0.75 | 0.86 | 1.0  | 0.71 | 0.92 | 0.64 | 0.9    | 0.64 | 0.56 | 0.71 | 0.86 | 0.86 | 0.78 | 0.73 |
| Ferry      | 1.0       | 0.57 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.86 | 1.0    | 0.57 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.86 |
| Floortile  | 0.52      | 0.68 | 0.64 | 0.82 | 0.83 | 0.91 | 0.66 | 0.80 | 0.68   | 0.68 | 0.89 | 0.73 | 1.0  | 0.82 | 0.86 | 0.74 |
| Grid       | 0.62      | 0.47 | 0.75 | 0.86 | 0.78 | 1.0  | 0.71 | 0.78 | 0.79   | 0.65 | 1.0  | 0.86 | 0.88 | 1.0  | 0.89 | 0.83 |
| Gripper    | 1.0       | 0.67 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.89 | 1.0    | 0.67 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.89 |
| Hanoi      | 1.0       | 0.50 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.83 | 0.75   | 0.75 | 1.0  | 1.0  | 1.0  | 1.0  | 0.92 | 0.92 |
| Miconic    | 0.75      | 0.33 | 0.50 | 0.50 | 0.75 | 1.0  | 0.67 | 0.61 | 0.89   | 0.89 | 1.0  | 0.75 | 0.75 | 1.0  | 0.88 | 0.88 |
| Satellite  | 0.60      | 0.21 | 1.0  | 1.0  | 1.0  | 0.75 | 0.87 | 0.65 | 0.82   | 0.64 | 1.0  | 1.0  | 1.0  | 0.75 | 0.94 | 0.80 |
| Transport  | 1.0       | 0.40 | 1.0  | 1.0  | 1.0  | 0.80 | 1.0  | 0.73 | 1.0    | 0.70 | 0.83 | 1.0  | 1.0  | 0.80 | 0.94 | 0.83 |
| Visitall   | 1.0       | 0.50 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.83 | 1.0    | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  |
| Zenotravel | 1.0       | 0.36 | 1.0  | 1.0  | 1.0  | 0.71 | 1.0  | 0.69 | 1.0    | 0.64 | 0.88 | 1.0  | 1.0  | 0.71 | 0.96 | 0.79 |
|            | 0.88      | 0.50 | 0.88 | 0.92 | 0.95 | 0.91 | 0.90 | 0.78 | 0.90   | 0.74 | 0.93 | 0.92 | 0.96 | 0.91 | 0.93 | 0.86 |

Table 1: *Precision and recall* obtained learning from labeled plans without (left) and with (right) static predicates.

|            | No Static |            |        | Static |            |        |
|------------|-----------|------------|--------|--------|------------|--------|
|            | Total     | Preprocess | Length | Total  | Preprocess | Length |
| Blocks     | 0.04      | 0.00       | 72     | 0.03   | 0.00       | 72     |
| Driverlog  | 0.14      | 0.09       | 83     | 0.06   | 0.03       | 59     |
| Ferry      | 0.06      | 0.03       | 55     | 0.06   | 0.03       | 55     |
| Floortile  | 2.42      | 1.64       | 168    | 0.67   | 0.57       | 77     |
| Grid       | 4.82      | 4.75       | 88     | 3.39   | 3.35       | 72     |
| Gripper    | 0.03      | 0.01       | 43     | 0.01   | 0.00       | 43     |
| Hanoi      | 0.12      | 0.06       | 48     | 0.09   | 0.06       | 39     |
| Miconic    | 0.06      | 0.03       | 57     | 0.04   | 0.00       | 41     |
| Satellite  | 0.20      | 0.14       | 67     | 0.18   | 0.12       | 60     |
| Transport  | 0.59      | 0.53       | 61     | 0.39   | 0.35       | 48     |
| Visitall   | 0.21      | 0.15       | 40     | 0.17   | 0.15       | 36     |
| Zenotravel | 2.07      | 2.04       | 71     | 1.01   | 1.00       | 55     |

Table 2: Total planning time, preprocessing time and plan length learning from labeled plans without/with static predicates.

## Evaluation

This section evaluates the performance of our approach for learning STRIPS action models starting from different amounts of available input knowledge.

**Setup.** The domains used in the evaluation are IPC domains that satisfy the STRIPS requirement (Fox and Long 2003), taken from the PLANNING.DOMAINS repository (Muise 2016). We only use 5 learning examples for each domain and they are fixed for all the experiments so we can evaluate the impact of the input knowledge in the quality of the learned models. All experiments are run on an Intel Core i5 3.10 GHz x 4 with 4 GB of RAM.

**Reproducibility.** We make fully available the compilation source code, the evaluation scripts and the used benchmarks at this anonymous repository <https://github.com/anonsub/strips-learning> so any experimental data reported in the paper is fully reproducible.

**Planner.** The classical planner we use to solve the instances that result from our compilations is MADAGASCAR (Rintanen 2014). We use MADAGASCAR because its

ability to deal with planning instances populated with dead-ends. In addition, SAT-based planners can apply the actions for programming preconditions in a single planning step (in parallel) because these actions do not interact. Actions for programming action effects can also be applied in a single planning step reducing significantly the planning horizon.

**Metrics.** The quality of the learned models is quantified with the *precision* and *recall* metrics. These two metrics are frequently used in *pattern recognition*, *information retrieval* and *binary classification* and are more informative than simply counting the number of errors in the learned model or computing the *symmetric difference* between the learned and the reference model (Davis and Goadrich 2006).

Intuitively, precision gives a notion of *soundness* while recall gives a notion of the *completeness* of the learned models. Formally,  $Precision = \frac{tp}{tp+fp}$ , where  $tp$  is the number of true positives (predicates that correctly appear in the action model) and  $fp$  is the number of false positives (predicates appear in the learned action model that should not appear). Recall is formally defined as  $Recall = \frac{tp}{tp+fn}$  where  $fn$  is the number of false negatives (predicates that should appear in the learned action model but are missing).

## Learning from labeled plans

We start evaluating our approach at  $\Lambda' = \langle \Psi, \Sigma, \Pi \rangle$  learning tasks, where *labeled plans* are available, and repeat the evaluation but exploiting potential static predicates that are computed from  $\Sigma$ . The potential *static predicates* considered here is the set of predicates s.t. every predicate instantiation appears unaltered in the initial and final states, for every  $\sigma_t \in \Sigma$ . Static predicates are used to constrain the space of possible action models as explained in the previous section.

Table 1 summarizes the obtained results. Precision (**P**) and recall (**R**) are computed separately for the preconditions (**Pre**), positive effects (**Add**) and negative Effects (**Del**) while the last two columns of each setting (and the last row) report averages values. We can observe that identifying static predicates drives to models with better precondition *recall*. This fact evidences that many of the missing precon-

ditions corresponded to static predicates because there were no incentive to learn them as they always hold (Gregory and Cresswell 2015).

Table 2 reports the total planning time, the preprocessing time (in seconds) invested by MADAGASCAR to solve the classical planning instances that result from our compilation as well as the number of actions in the solutions. All the learning tasks are solved in a few seconds time and interestingly, one can identify the domains with static predicates by just looking at the reported plan length. In these domains some preconditions corresponding to static predicates are directly derived from the learning examples so less programming actions are required. When static predicates are identified, the resulting compilation is also much compact and produces smaller planning/instantiation times.

### Learning from partially specified action models

We evaluate now the ability of our approach to support partially specified action models; that is, when addressing learning tasks of the kind  $\Lambda'' = \langle \Psi, \Sigma, \Pi, \Xi_0 \rangle$ . In this experiment, the model of half of the actions is given in  $\Xi_0$  as an extra input of the learning task.

Tables 3 and 4 summarize the obtained results, including the identification of static predicates. We report only the precision and recall for the *unknown* actions (considering the *known* action models would mean reporting higher scores since *precision* and *recall* in these models is 1.0). We must also note that errors in  $\Lambda''$  tasks have greater impact than in the corresponding  $\Lambda'$  tasks because the evaluation is just done over half of the actions (e.g. the preconditions recall for the *Floortile* or *Gripper* domains).

Remarkably, the overall precision is now 0.98 which means that the content of the learned models is highly reliable. Recall is a lower, 0.87, indicating that the learned models still miss some information. Particularly, preconditions are again the component more difficult to be fully learned. Overall the reported results confirm the previous trend: more input knowledge drives to better models and requires less planning time. Likewise smaller solution plans are required since it is only necessary to program half of the actions (the other half is input knowledge). *Visitall* and *Hanoi* are excluded from this evaluation because they only contain one action schema.

### Learning from (initial,final) state pairs

Finally we evaluate our approach at learning tasks where input plans are not available, so the planner must not only compute the action models but also the plans that satisfy the input labels. Table 5 and 6 summarize the obtained results using static predicates and partially specified action models so the inputs to this learning task are  $\Psi, \Sigma$  and  $\Xi_0$ .

Values for the *Zenotrail* and *Grid* domains are not reported because MADAGASCAR could not solved the corresponding planning tasks within a 1000 secs time bound. Precision and recall have now low values. The learning hypothesis space is now so low constrained that actions can be reformulated and still be compliant with the inputs (e.g. the *blocksworld* operator *stack* could be *learned* with the preconditions and effects of the *unstack* operator and vice

|           | Pre  |      | Add |      | Del |      |      |      |
|-----------|------|------|-----|------|-----|------|------|------|
|           | P    | R    | P   | R    | P   | R    | P    | R    |
| Blocks    | 1.0  | 1.0  | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 1.0  |
| Driverlog | 1.0  | 0.71 | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 0.90 |
| Ferry     | 1.0  | 0.67 | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 0.89 |
| Floortile | 0.75 | 0.60 | 1.0 | 0.80 | 1.0 | 0.80 | 0.92 | 0.73 |
| Grid      | 1.0  | 0.67 | 1.0 | 1.0  | 1.0 | 1.0  | 0.84 | 0.78 |
| Gripper   | 1.0  | 0.50 | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 0.83 |
| Miconic   | 1.0  | 1.0  | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 1.0  |
| Satellite | 1.0  | 0.57 | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 0.86 |
| Transport | 1.0  | 0.75 | 1.0 | 1.0  | 1.0 | 1.0  | 1.0  | 0.92 |
| Zenotrail | 1.0  | 0.67 | 1.0 | 1.0  | 1.0 | 0.67 | 1.0  | 0.78 |
|           | 0.98 | 0.71 | 1.0 | 0.98 | 1.0 | 0.95 | 0.98 | 0.87 |

Table 3: *Precision* and *recall* learning from partially specified action models.

|           | Total time | Preprocess | Plan length |
|-----------|------------|------------|-------------|
| Blocks    | 0.07       | 0.01       | 54          |
| Driverlog | 0.03       | 0.01       | 40          |
| Ferry     | 0.06       | 0.03       | 45          |
| Floortile | 0.43       | 0.42       | 55          |
| Grid      | 3.12       | 3.07       | 53          |
| Gripper   | 0.03       | 0.01       | 35          |
| Miconic   | 0.03       | 0.01       | 34          |
| Satellite | 0.14       | 0.14       | 47          |
| Transport | 0.23       | 0.21       | 37          |
| Zenotrail | 0.90       | 0.89       | 40          |

Table 4: Time and plan length learning from partially specified action models.

versa). We tried to minimize this effect with the additional input knowledge (static predicates and partially specified action models) but still the achieved scores are below the scores obtained when learning from labeled plans (despite no other input knowledge was provided).

### Past and future of learning action models

Back in the 90's various systems aimed learning operators mostly via interaction with the environment. LIVE captured and formulated observable features of objects and used them to acquire and refine operators (Shen and Simon 1989). OBSERVER updated preconditions and effects by removing and adding facts, respectively, accordingly to observations (Wang 1995). These early works were based on lifting the observed states supported by exploratory plans or external teachers.

Action model learning has also been studied in domains where there is partial or missing state observability. ARMS works when no partial intermediate state is given. It defines a set of weighted constraints that must hold for the plans to be correct, and solves the weighted propositional satisfiability problem with a MAX-SAT solver (Yang, Wu, and Jiang 2007). In order to efficiently solve the large MAX-SAT representations, ARMS implements a hill-climbing method that models the actions approximately. SLAF also deals with partial observability (Amir and Chang 2008). Given a formula representing the initial belief state, a sequence of executed actions and the corresponding partially observed

states, it builds a complete explanation of observations by models of actions through a CNF formula. The learning algorithm updates the formula of the belief state with every action and observation in the sequence. This update makes sure that the new formula represents all the transition relations consistent with the actions and observations. The formula returned at the end includes all consistent models, which can then be retrieved with additional processing.

LOCM only requires the example plans as input without need for providing information about predicates or states (Cresswell, McCluskey, and West 2013). This makes LOCM be most likely the learning approach that works with the least information possible. The lack of available information is addressed by LOCM by exploiting assumptions about the kind of domain model it has to generate. Particularly, it assumes a domain consists of a collection of objects (sorts) whose defined set of states can be captured by a parameterized Finite State Machine (FSM). The intuitive assumptions of LOCM yield a learning model heavily reliant on the kind of domain structure. The inability of LOCM to properly derive domain theories where the state of a sort is subject to different FSMs is later overcome by LOCM2 by forming separate FSMs, each containing a subset of the full transition set for the sort (Cresswell and Gregory 2011). LOP (LOCM with Optimized Plans (Gregory and Cresswell 2015)), the last contribution of the LOCM family, addresses the problem of inducing static predicates. Because LOCM approaches induce similar models for domains with similar structures, they face problems at generating models for domains that are only distinguished by whether or not they contain static relations (e.g. *blocksworld* and *freecell*). In order to mitigate this drawback, LOP applies a post-processing step after the LOCM analysis which requires additional information about the plans, namely a set of optimal plans to be used in the learning phase.

Compiling the learning of action models into classical planning is a general and flexible approach that allows to accommodate different amounts and kinds of input knowledge and opens up a path for addressing further learning and validation tasks. For instance, the example plans in  $\Pi$  could be replaced or complemented by a set  $\mathcal{O}$  of sequences of observations (i.e., fully or partial state observations with noisy or missing fluents (Sohrabi, Riabov, and Udrea 2016)), so learning tasks  $\Lambda = \langle \Psi, \Sigma, \mathcal{O}, \Xi_0 \rangle$  could also be addressed. Furthermore, our approach seems extensible to learning other types of generative models (e.g. hierarchical models like HTN or behaviour trees), that can be more appealing than STRIPS models, since using them to compute planning solutions requires less search effort.

## Conclusions

We presented a novel approach for learning STRIPS action models from examples using classical planning. The approach is flexible to different amounts of available input knowledge and accepts partially specified action models. We also introduced *precision* and *recall* metrics, widely used in ML, to the evaluation of planning action models. These metrics separately assess the correctness and completeness of the learned models making easier the identification of flaws.

|            | Pre  |      | Add  |      | Del  |      |      |      |
|------------|------|------|------|------|------|------|------|------|
|            | P    | R    | P    | R    | P    | R    | P    | R    |
| Blocks     | 0.33 | 0.33 | 0.75 | 0.50 | 0.33 | 0.33 | 0.47 | 0.39 |
| Driverlog  | 1.0  | 0.29 | 0.33 | 0.67 | 1.0  | 0.50 | 0.78 | 0.48 |
| Ferry      | 1.0  | 0.67 | 0.50 | 1.0  | 1.0  | 1.0  | 0.83 | 0.89 |
| Floortile  | 0.67 | 0.40 | 0.50 | 0.40 | 1.0  | 0.40 | 0.72 | 0.40 |
| Grid       | -    | -    | -    | -    | -    | -    | -    | -    |
| Gripper    | 1.0  | 0.50 | 1.0  | 1.0  | 1.0  | 1.0  | 1.0  | 0.83 |
| Miconic    | 0.0  | 0.0  | 0.33 | 0.50 | 0.0  | 0.0  | 0.11 | 0.17 |
| Satellite  | 1.0  | 0.14 | 0.67 | 1.0  | 1.0  | 1.0  | 0.89 | 0.71 |
| Transport  | 0.0  | 0.0  | 0.25 | 0.5  | 0.0  | 0.0  | 0.08 | 0.17 |
| Zenotravel | -    | -    | -    | -    | -    | -    | -    | -    |
|            | 0.63 | 0.29 | 0.54 | 0.70 | 0.67 | 0.53 | 0.61 | 0.51 |

Table 5: *Precision* and *recall* learning from state pairs.

|            | Total time | Preprocess | Plan length |
|------------|------------|------------|-------------|
| Blocks     | 2.14       | 0.00       | 58          |
| Driverlog  | 0.09       | 0.00       | 88          |
| Ferry      | 0.17       | 0.01       | 65          |
| Floortile  | 6.42       | 0.15       | 126         |
| Grid       | -          | -          | -           |
| Gripper    | 0.03       | 0.00       | 47          |
| Miconic    | 0.04       | 0.00       | 68          |
| Satellite  | 4.34       | 0.10       | 126         |
| Transport  | 2.57       | 0.21       | 47          |
| Zenotravel | -          | -          | -           |

Table 6: Time and plan length when learning from state pairs.

As far as we know, this is the first work on learning action models exclusively using an *off-the-shelf* classical planner and evaluated over a wide range of different domains. Recently, ?? proposed a classical planning compilation for learning action models but following the *finite domain* representation for the state variables and did not report experimental results since the compilation was not implemented.

Our evaluation shows that, when example plans are available, we can compute accurate action models from small sets of learning examples (only five learning examples per domain) and investing small learning times (less than a second time in most of the domains). When action plans are not available, our approach is still able to produce action models compliant with the input information. In this case, since learning is not constrained by actions it can reformulate operators changing their semantics and making tricky their comparison with a reference model.

The size of the compiled classical planning instances depend on the number of input examples. Generating *informative* examples for learning planning action models is an open issue. Planning actions include preconditions that are only satisfied by specific sequences of actions, often, with a low probability of being chosen by chance (Fern, Yoon, and Givan 2004). The success of recent algorithms for exploring planning tasks (?) motivates the development of novel techniques able to autonomously collect informative learning examples. The combination of such exploration techniques with our learning approach is an intriguing research direction that opens the door to the bootstrapping of planning action models.



## Acknowledgment

This work is supported by the Spanish MINECO project TIN2017-88476-C2-1-R. Diego Aineto is partially supported by the FPU16/03184 and Sergio Jiménez by the *Ramon y Cajal*, both programs funded by the Spanish government.

## References

- Amir, E., and Chang, A. 2008. Learning partially observable deterministic action models. *Journal of Artificial Intelligence Research* 33:349–402.
- Bonet, B.; Palacios, H.; and Geffner, H. 2009. Automatic derivation of memoryless policies and finite-state controllers using classical planners. In *International Conference on Automated Planning and Scheduling (ICAPS)*.
- Cresswell, S., and Gregory, P. 2011. Generalised domain model acquisition from action traces. In *International Conference on Automated Planning and Scheduling, ICAPS-11*.
- Cresswell, S. N.; McCluskey, T. L.; and West, M. M. 2013. Acquiring planning domain models using LOCM. *The Knowledge Engineering Review* 28(02):195–213.
- Davis, J., and Goadrich, M. 2006. The relationship between precision-recall and ROC curves. In *International Conference on Machine Learning*, 233–240. ACM.
- Fern, A.; Yoon, S. W.; and Givan, R. 2004. Learning domain-specific control knowledge from random walks. In *International Conference on Automated Planning and Scheduling, ICAPS-04*, 191–199.
- Fikes, R. E., and Nilsson, N. J. 1971. Strips: A new approach to the application of theorem proving to problem solving. *Artificial Intelligence* 2(3-4):189–208.
- Fox, M., and Long, D. 1998. The automatic inference of state invariants in TIM. *Journal of Artificial Intelligence Research* 9:367–421.
- Fox, M., and Long, D. 2003. PDDL2.1: An extension to PDDL for expressing temporal planning domains. *Journal of Artificial Intelligence Research* 20:61–124.
- Geffner, H., and Bonet, B. 2013. A concise introduction to models and methods for automated planning.
- Ghallab, M.; Nau, D.; and Traverso, P. 2004. *Automated Planning: theory and practice*. Elsevier.
- Gregory, P., and Cresswell, S. 2015. Domain model acquisition in the presence of static relations in the LOP system. In *International Conference on Automated Planning and Scheduling, ICAPS-15*, 97–105.
- Howey, R.; Long, D.; and Fox, M. 2004. VAL: Automatic plan validation, continuous effects and mixed initiative planning using PDDL. In *Tools with Artificial Intelligence, 2004. ICTAI 2004. 16th IEEE International Conference on*, 294–301. IEEE.
- Kambhampati, S. 2007. Model-lite planning for the web age masses: The challenges of planning with incomplete and evolving domain models. In *National Conference on Artificial Intelligence, AAAI-07*.
- McDermott, D.; Ghallab, M.; Howe, A.; Knoblock, C.; Ram, A.; Veloso, M.; Weld, D.; and Wilkins, D. 1998. PDDL – The Planning Domain Definition Language.
- Michalski, R. S.; Carbonell, J. G.; and Mitchell, T. M. 2013. *Machine learning: An artificial intelligence approach*. Springer Science & Business Media.
- Muise, C. 2016. Planning. domains. *ICAPS system demonstration*.
- Ramírez, M. 2012. *Plan recognition as planning*. Ph.D. Dissertation, Universitat Pompeu Fabra.
- Rintanen, J. 2014. Madagascar: Scalable planning with sat. *Proceedings of the 8th International Planning Competition (IPC-2014)*.
- Segovia-Aguas, J.; Jiménez, S.; and Jonsson, A. 2016. Hierarchical finite state controllers for generalized planning. In *International Joint Conference on Artificial Intelligence, IJCAI-16*, 3235–3241. AAAI Press.
- Segovia-Aguas, J.; Jiménez, S.; and Jonsson, A. 2017. Generating context-free grammars using classical planning. In *International Joint Conference on Artificial Intelligence, ICAPS-17*.
- Shen, W., and Simon, H. A. 1989. Rule creation and rule learning through environmental exploration. In *International Joint Conference on Artificial Intelligence, IJCAI-89*, 675–680.
- Slaney, J., and Thiébaux, S. 2001. Blocks world revisited. *Artificial Intelligence* 125(1-2):119–153.
- Sohrabi, S.; Riabov, A. V.; and Udrea, O. 2016. Plan recognition as planning revisited. In *International Joint Conference on Artificial Intelligence, IJCAI-16*, 3258–3264.
- Vallati, M.; Chrapa, L.; Grzes, M.; McCluskey, T. L.; Roberts, M.; and Sanner, S. 2015. The 2014 international planning competition: Progress and trends. *AI Magazine* 36(3):90–98.
- Wang, X. 1995. Learning by observation and practice: An incremental approach for planning operator acquisition. In *International Conference on Machine Learning*, 549–557.
- Yang, Q.; Wu, K.; and Jiang, Y. 2007. Learning action models from plan examples using weighted max-sat. *Artificial Intelligence* 171(2-3):107–143.