

One-Shot Learning of Temporal Actions Models via Constraint Programming

Antonio Garrido and Sergio Jiménez

Abstract. We present a *Constraint Programming* (CP) formulation for learning temporal planning action models from the observation of a single plan execution (*one shot*). Inspired by the CSP approach to *temporal planning*, our CP formulation models *time-stamps* for states and actions, *causal-link* relationships, *threats* and effect *interferences*. This modeling evidences the connection between the tasks of *plan synthesis*, *plan validation* and *action model learning* in the temporal planning setting. The CP formulation is solver-independent so off-the-shelf CSP solvers can be used for the resolution of any of these three tasks. The performance of our CP formulation is assessed when *learning* and *validating* action models at several temporal planning domains specified in the PDDL2.1 representation language.

1 INTRODUCTION

Temporal planning is an expressive planning model that relaxes the assumption of instantaneous actions of *classical planning* [10]. Actions in temporal planning are called *durative*, have an associated duration and hence, their conditions/effects may hold/happen at different times [7]. This means that *durative actions* can be executed in parallel and overlap in several different ways [4], and that valid solutions for temporal planning instances specify the precise time-stamp when each durative action start and end [16].

Despite the potential of state-of-the-art planners, their application to real world problems is still somewhat limited mainly because of the difficulty of specifying correct and complete planning models [18]. The more expressive the planning model, the more evident becomes the *knowledge acquisition bottleneck* that jeopardizes the usability of planning technology. There are growing efforts in the planning community for the machine learning of action models from sequential plans: since pioneering learning systems like ARMS [23], we have seen systems able to learn action models with *quantifiers* [2, 28], from *noisy* actions or states [20, 24], from *null state information* [3], or from *incomplete* domain models [25, 27].

Most of the cited approaches for model learning are purely inductive and require large input datasets, e.g. hundreds of plan observations, to compute statistically significant models and focus on learning models from sequential plans for classical planning. With the aim of understanding better the connection between the learning of durative action models, *temporal planning* and the validation of temporal plans, this paper follows a radically different approach and studies the singular learning scenario where just the observation of a single plan execution (*one-shot*) is available. We leverage a solver-independent CP formulation that connects the *learning* of planning action models with the *synthesis* and the *validation* of temporal plans. Off-the-shelf CSP solvers can be used for any of these tasks.

As far as we know this paper presents the first approach for learn-

ing action models for the *temporal planning* setting. While learning an action model for classical planning means computing the actions' conditions and effects that are consistent with the input observations, learning temporal action models requires additionally: i) identifying how conditions and effects are temporally distributed within the actions, and ii) estimate the action duration. Further our approach allows the learning of action models from observations of plans with overlapping actions. This feature makes our approach appealing for learning action models from observations of multi-agent environments [8].

2 BACKGROUND

This section formalizes the *temporal planning* and the *Constraint Satisfaction* model that we follow in this work.

2.1 Temporal Planning

We assume that *states* are factored into a set F of Boolean variables. A state s is a time-stamped assignment of values to all the variables in F . A *temporal planning problem* is a tuple $P = \langle F, I, G, A \rangle$ where the *initial state* I is a fully observed state (i.e. a total assignment of the state variables $|I| = |F|$) time-stamped with $t = 0$; $G \subseteq F$ is a conjunction of *goal conditions* over the variables in F that defines the set of goal states; and A represents the set of *durative actions*.

A *durative action* has an associated duration and may have conditions/effects on F at different times [9, 22]. To compactly represent temporal planning problems, we assume that the state variables in F are instantiations of a given set of predicates Ψ (like in the PDDL language) and that durative actions in A are fully grounded from *action schemes* (also known as *operators*).

PDDL2.1 is the input representation language for the temporal track of the International Planning Competition (IPC) [7, 12]. According to PDDL2.1, a durative action $a \in A$ is defined with the following elements:

1. $\text{dur}(a)$, a positive value indicating the *duration* of the action.
2. $\text{cond}_s(a)$, $\text{cond}_o(a)$, $\text{cond}_e(a)$ representing the three types of action *conditions*. Unlike the *preconditions* of classical actions, action conditions in PDDL2.1 must hold: before a is executed (*at start*), during the entire execution of a (*over all*) or when a finishes (*at end*), respectively.
3. $\text{eff}_s(a)$ and $\text{eff}_e(a)$ represent the two types of action effects. In PDDL2.1, effects can happen *at start* or *at end* of action a respectively, and can be either positive or negative (i.e. asserting or retracting variables).

PDDL2.1 is a restricted temporal planning model that defines the semantics of a *durative action* a as two discrete events, $\text{start}(a)$ and

$\text{end}(a) = \text{start}(a) + \text{dur}(a)$. This means that if a starts on state s with time-stamp $\text{start}(a)$, then $\text{cond}_s(a)$ must hold in s . Ending action a in state s' , with time-stamp $\text{end}(a)$, means $\text{cond}_e(a)$ must hold in s' . Over all conditions must hold at any state between s and s' or, in other words, throughout the closed interval $[\text{start}(a), \text{end}(a)]$. Likewise, *at start* and *at end* effects are instantaneously applied at states s and s' , respectively (continuous effects are not considered in this work). Figure 1 shows an example of two schemes for PDDL2.1 durative actions taken from the *driverlog* domain. The schema `board-truck` has fixed duration of two time units while the duration of `drive-truck` depends on the driving time associated to the two given locations.

```
(:durative-action board-truck
:parameters (?d - driver ?t - truck ?l - location)
:duration (= ?duration 2)
:condition (and (at start (at ?d ?l))
                (at start (empty ?t))
                (over all (at ?t ?l)))
:effect (and (at start (not (at ?d ?l)))
             (at start (not (empty ?t)))
             (at end (driving ?d ?t))))

(:durative-action drive-truck
:parameters (?t - truck ?l1 - location ?l2 - location
             ?d - driver)
:duration (= ?duration (driving-time ?l1 ?l2))
:condition (and (at start (at ?t ?l1))
                (at start (link ?l1 ?l2))
                (over all (driving ?d ?t)))
:effect (and (at start (not (at ?t ?l1)))
             (at end (at ?t ?l2))))
```

Figure 1. Two action schemes of durative actions represented in PDDL2.1.

PDDL2.2 is an extension of the PDDL2.1 language that includes the notion of *Timed Initial Literal* [14], denoted as $\text{til}(f, t)$, and representing that variable $f \in F$ becomes true at a certain time $t > 0$, independently of the actions in the plan [5]. TILs are useful to model *exogenous events*; for instance, in a logistics scenario, the time window when a warehouse is open can be modeled with these two timed initial literals: $\text{til}(\text{open}, 8)$ and $\text{til}(\neg \text{open}, 20)$.

A *temporal plan* is a set of pairs $\pi = \{(a_1, t_1), (a_2, t_2) \dots (a_n, t_n)\}$. Each pair (a_i, t_i) contains a *durative action* a_i and the action *time-stamp* $t_i = \text{start}(a_i)$. The execution of a temporal plan starting from a given initial state I induces a state sequence formed by the union of all states $\{s_{t_i}, s_{t_i + \text{dur}(a_i)}\}$, where there exists an initial state $s_0 = I$, and a state s_{end} that is the last state induced by the execution of the plan. Sequential plans can then be expressed as temporal plans but not the opposite. A *solution* to a given temporal planning problem P is a *temporal plan* π such that its execution, starting from the corresponding initial state, eventually reaches a state that meets the goal conditions, $G \subseteq s_{\text{end}}$. A solution is *optimal* iff it minimizes the plan *makespan* (i.e., the maximum $\text{end}(a) = \text{start}(a) + \text{dur}(a)$ of any actions in the plan).

2.2 Constraint Satisfaction

A *Constraint Satisfaction Problem* (CSP) is a tuple $\langle X, D, C \rangle$, where X is a set of finite-domain *variables*, D represents the *domain* for each of these variables and C is a set of *constraints* among the variables in X that bound their possible values in D .

A *solution* to a CSP is an assignment of values to all the variables in X that is *consistent* with all the input constraints. Given a CSP there may be many different solutions to that problem, i.e., different variable assignments that are *consistent* with the input constraints.

A *cost-function* can be defined over the variables in X to specify user preferences about the space of possible solutions. Given a CSP and a cost-function, then an *optimal solution* is a full variable assignment that is consistent with the constraints of the CSP while it minimizes the value of the defined cost-function.

3 One-shot learning of temporal actions models

We formalize the task of the *one-shot learning of temporal action models* as a tuple $\mathcal{L} = \langle F, I, G, A?, O, C \rangle$ where:

- $\langle F, I, G, A? \rangle$ is a *temporal planning problem* such that the actions in $A?$ are *partially specified*. This means that the exact conditions/effects, their temporal annotation, and the duration of actions are unknown while the actions *header* (i.e., the *name* and *parameters* of each action) is known. With this regard, we say that a fluent $f \in F$ is *candidate* to appear in the condition/effects of an action $a \in A?$ iff f appears in the set of FOL interpretations of the predicates that shape the fluents F over the action parameters $\text{pars}(a)$. For instance, Figure 2 shows the set of six *candidates* to appear in the conditions/effect of the ground action `board-truck(driver1, truck1, loc1)`.

```
;;;
;; Candidates for board-truck(driver1, truck1, loc1)

(at driver1 loc1) (at truck1 loc1)
(driving driver1 truck1) (empty truck1)
(path loc1 loc1) (link loc1 loc1)
```

Figure 2. Set of six *candidates* to appear in the conditions/effects of the ground action `board-truck(driver1, truck1, loc1)`.

- O is the set of *observations* over a single plan execution. At least this set contains a full observation of the initial state (time-stamped with $t = 0$) and a final state observation, that equals the goals G of the temporal planning problem, time-stamped with t_{end} (the *makespan* of the observed plan). Additionally, it can contain time-stamped observations of traversed intermediate *partial states*¹ as well as the times when actions start and/or end their execution. For instance, Figure 3 shows an example of the observation of a plan execution taken from the *driverlog* domain.
- C is a set of *constraints* that captures domain-specific expert knowledge. In this work these constraints are of two kinds:
 - Constraints that specify that a given *candidate* $f \in F$ is actually in the conditions/effects of a given action. These constraints allow to represent partially specified action models [26]. For instance we may know in advance that the action `board-truck` requires the `driver` and the `truck` to be at the same location.
 - Mutually-exclusive (*mutex*) constraints that allow to (1), deduce new observations and (2), prune action models inconsis-

¹ In this work, not all variables can be observed at any time; that is, we deal with *partial observations* (e.g. just a subset of variables is observable by associated sensors). Observations are noiseless, which means that if a value is observed, that is the actual value of that variable.

```

(objects driver1 driver2 - driver
 truck1 truck2 - truck
 package1 package2 - obj
 s0 s1 s2 p1-0 p1-2 - location)

(:init (at driver1 s2) (at driver2 s2) (at truck1 s0)
 (empty truck1) (at truck2 s0) (empty truck2)
 (at package1 s0) (at package2 s0)
 (path s1 p1-0) (path p1-0 s1) (path s0 p1-0)
 (path p1-0 s0) (path s1 p1-2) (path p1-2 s1)
 (path s2 p1-2) (path p1-2 s2)
 (link s0 s1) (link s1 s0) (link s0 s2) (link s2 s0)
 (link s2 s1) (link s1 s2))

(:observation :time-stamp 56
 (at driver1 s1) (at truck1 s1))

(:observation :time-stamp 78
 (at package1 s0) (at package2 s0))

```

Figure 3. Example of a set of three observations (containing the fully observed initial state and two time-stamped partial states) extracted from the execution of a plan from the *driverlog* domain.

tent with these constraints. Figure 4 shows an example of a set of five mutex constraints for the *driverlog* domain.

A *solution* for the learning task \mathcal{L} is a fully specified model of durative actions \mathcal{A} such that the *conditions*, *effects*, their temporal annotations and the *duration* of any action in \mathcal{A} are: i) completely specified; and ii) *consistent* with $\mathcal{L} = \langle F, I, G, A?, O, C \rangle$. By *consistent* we mean that there exists a valid plan that exclusively contains actions in \mathcal{A} and whose execution starting in I , produces all the observations in O at the associated time-stamps, while it satisfies all constraints in C , and reaches a final state that satisfies G .

```

 $\forall \text{truck}, \text{driver} : \neg \text{empty}(\text{truck}) \vee \neg \text{driving}(\text{driver}, \text{truck}).$ 
 $\forall \text{driver}, \text{loc}_1, \text{loc}_2 : \neg \text{at}(\text{driver}, \text{loc}_1) \vee \neg \text{at}(\text{driver}, \text{loc}_2),$ 
 $\neq (\text{loc}_1, \text{loc}_2).$ 
 $\forall \text{driver}, \text{truck}_1, \text{truck}_2 : \neg \text{driving}(\text{driver}, \text{truck}_1) \vee$ 
 $\neg \text{driving}(\text{driver}, \text{truck}_2), \neq (\text{truck}_1, \text{truck}_2).$ 
 $\forall \text{drv}_1, \text{drv}_2, \text{truck} : \neg \text{driving}(\text{drv}_1, \text{truck}) \vee$ 
 $\neg \text{driving}(\text{drv}_2, \text{truck}), \neq (\text{drv}_1, \text{drv}_2).$ 
 $\forall \text{drv}, \text{location}, \text{truck} : \neg \text{at}(\text{drv}, \text{location}) \vee \neg \text{driving}(\text{drv}, \text{truck}).$ 

```

Figure 4. Examples of five *mutex constraints* for the *driverlog* domain.

4 One-Shot learning of action models with CSPs

Given a *one-shot learning task* \mathcal{L} as defined in Section 3, we automatically create a CSP, whose solution induces an action model that solves \mathcal{L} . This method is solver-independent and integrates previous work on *temporal planning* as satisfiability [22, 9, 21].

4.1 The CSP variables

For each action $a \in A?$ and *candidate* $f \in F$ to appear in the conditions/effects of a , we create the following eight CSP variables that are shown in Table 1.

Variable X1 represents the time when an action *starts* (its time-stamp), X2 represents when the action *ends* and variable X3 represents the action duration. The value of X1, X2 and X3 can be either observed in O or derived from the expression $\text{end}(a) = \text{start}(a) + \text{dur}(a)$. We model time in \mathbb{Z}^+ and bound all maximum times to the *makespan* of the observed plan (t_{end} if observed in O). If the observation of t_{end} is unavailable, we consider a large enough domain for time. Boolean variables X4/X5 model whether f is actually a condition/effect of action a . X6.1 and X6.2 define the closed

interval throughout condition f must hold for the application of action a (provided that $\text{is_cond}(f, a) = \text{true}$). X7 models a *causal link* representing that action b supports f that is required by a . If f is not a condition of a ($\text{is_cond}(f, a) = \text{false}$) then $\text{sup}(f, a) = \emptyset$, thus representing an empty supporter. Last but not least variable X8 models the time-stamp when effect f happens in a (provided $\text{is_eff}(f, a) = \text{true}$).

Table 1. The CSP variables, their domains and semantics.

ID	Variable	Domain	Description
X1	$\text{start}(a)$	$[0..t_{\text{end}}]$	Start time of action a
X2	$\text{end}(a)$	$[0..t_{\text{end}}]$	End time of action a
X3	$\text{dur}(a)$	$[0..t_{\text{end}}]$	Duration of action a
X4	$\text{is_cond}(f, a)$	$\{0, 1\}$	1 if f is a condition of a ; 0 otherwise
X5	$\text{is_eff}(f, a)$	$\{0, 1\}$	1 if f is an effect of a ; 0 otherwise
X6.1	$\text{req_start}(f, a)$	$[0..t_{\text{end}}]$	Interval when action a requires f
X6.2	$\text{req_end}(f, a)$		
X7	$\text{sup}(f, a)$	$\{b \mid b \in A? \cup \emptyset\}$	Supporters for causal link $\langle b, f, a \rangle$
X8	$\text{time}(f, a)$	$[0..t_{\text{end}}]$	Time when the effect f of a happens

This simple formulation is able to model *tils* and *observations*. The intuition is that modeling a *til* is analogous to modeling the *initial state* of a planning task (both represent information that is given at a particular time but externally to the execution of the plan). Likewise modeling an observation is analogous to modeling the goal of a planning task, as they both represent conditions that must be satisfied in the execution of the plan. On the one hand a *til*(f, t) is modeled as a *dummy* action that starts at time t and has instantaneous duration ($\text{start}(\text{til}(f, t)) = t$ and $\text{dur}(\text{til}(f, t)) = 0$) with no conditions and the single effect f that happens at time t ($\text{is_eff}(f, \text{til}(f, t)) = \text{true}$ and $\text{time}(f, \text{til}(f, t)) = t$). On the other hand, an observation $\text{obs}(f, t)$ is modeled as another *dummy* action that also starts at time t and has instantaneous duration ($\text{start}(\text{obs}(f, t)) = t$ and $\text{dur}(\text{obs}(f, t)) = 0$) but with only one condition f , which is the value observed for fact f ($\text{is_cond}(f, \text{obs}(f, t)) = \text{true}$, $\text{sup}(f, \text{obs}(f, t)) \neq \emptyset$ and $\text{req_start}(f, \text{obs}(f, t)) = \text{req_end}(f, \text{obs}(f, t)) = t$), and no effects at all. Observations can also refer to the start and end of an action, $\text{obs}(\text{is_start}(a), t)$ represents that action a starts at t while $\text{obs}(\text{is_end}(a), t)$ represents that action a ends at t .

4.2 The CSP constraints

Table 2 shows the constraints defined among the CSP variables of Table 1. C1 models the duration of an action while C2 indicates that actions must end before t_{end} . C3 forces to have a well-defined interval $[\text{req_start}, \text{req_end}]$, when the conditions of action a are required. C4 models that only action conditions require supporters and C5 models that the time when b supports f must be before a requires it because of the causal link $\langle b, f, a \rangle^2$. Given a causal link $\langle b, f, a \rangle$, constraint C6 avoids *threats* of actions c deleting f (threats are solved via *promotion* or *demotion* [12]). C7 prevents action a being a supporter of f when $\text{is_eff}(f, a) = \text{false}$. Constraint C8 models the fact that when the same action requires and deletes f the effect cannot happen before the condition. Note the \geq inequality here: if one condition and one effect of the same action happen at the same time, the underlying semantics in planning considers the condition is checked instantly before the effect [7]. C9 prevents two actions have contradictory effects and C10 forces actions to have at least one condition and one effect (C9 applies to any type of action, including the dummy actions *init*, *goal*, *til* and *obs*, while C10 only applies to *non-dummy* actions).

² $\text{time}(f, b) < \text{req_start}(f, a)$ and not \leq because our temporal planning model assumes $\epsilon > 0$ (ϵ denotes a small tolerance that implies no collision between the time when effect f is supported and when it is required, like in PDDL2.1 [7]). When time is modeled in \mathbb{Z}^+ , $\epsilon = 1$ so \leq becomes $<$.

Table 2. The CSP constraints and a brief description.

ID	Constraint	Description
C1	$\text{end}(a) = \text{start}(a) + \text{dur}(a)$	Relationship among start, end and duration of a
C2	$\text{end}(a) \leq \text{start}(\text{goal})$	Always goal is the last action of the plan
C3	if $(\text{is_cond}(f, a) = \text{true})$ then $\text{req_start}(f, a) \leq \text{req_end}(f, a)$	$[\text{req_start}(f, a) \dots \text{req_end}(f, a)]$ is a valid interval
C4	iff $(\text{is_cond}(f, a) = \text{false})$ then $\text{sup}(f, a) = \emptyset$	f is not a condition of $a \iff$ the supporter of f in a is \emptyset
C5	if $(\text{is_eff}(f, b) = \text{true})$ AND $(\text{is_cond}(f, a) = \text{true})$ AND $(\text{sup}(f, a) = b)$ then $\text{time}(f, b) < \text{req_start}(f, a)$	Modeling the causal link $\langle b, f, a \rangle$: supporting f before it is required (obviously $b \neq \emptyset$)
C6	if $(\text{is_eff}(f, b) = \text{true})$ AND $(\text{is_cond}(f, a) = \text{true})$ AND $(\text{is_eff}(\text{not-}f, c) = \text{true})$ AND $(\text{sup}(f, a) = b)$ AND $(c \neq a)$ then $(\text{time}(\text{not-}f, c) < \text{time}(f, b))$ OR $(\text{time}(\text{not-}f, c) > \text{req_end}(f, a))$	Solving threat of c to causal link $\langle b, f, a \rangle$ by promotion or demotion (obviously $b \neq \emptyset$)
C7	if $(\text{is_eff}(f, a) = \text{false})$ then $\forall b$ that requires f : $\text{sup}(f, b) \neq a$	a cannot be a supporter of f for any other action b
C8	if $(\text{is_cond}(f, a) = \text{true})$ AND $(\text{is_eff}(\text{not-}f, a) = \text{true})$ then $\text{time}(\text{not-}f, a) \geq \text{req_end}(f, a)$	a requires and deletes f : the condition holds before the effect
C9	if $(\text{is_eff}(f, b) = \text{true})$ AND $(\text{is_eff}(\text{not-}f, c) = \text{true})$ then $\text{time}(f, b) \neq \text{time}(\text{not-}f, c)$	Solving effect interference at the same time (f and $\text{not-}f$)
C10	$\sum \text{is_cond}(f_i, a) \geq 1$ AND $\sum \text{is_eff}(f_j, a) \geq 1$ forall condition f_i and effect f_j of a	Every non-dummy action has at least one condition/effect

4.2.1 Constraints for the PDDL2.1 model

The presented CSP formulation accommodates a level of expressiveness beyond PDDL2.1 because it allows conditions/effects to be at any time, even outside the execution of the action. For example, it allows a condition f to hold in $\text{start}(a) \pm 2$: $\text{req_start}(f, a) = \text{start}(a) - 2$ and $\text{req_end}(f, a) = \text{start}(a) + 2$. Likewise an effect f might also happen after the action ends e.g., $\text{time}(f, a) = \text{end}(a) + 2$.

Making the formulation PDDL2.1-compliant is straightforward, by adding the constraints of Table 3 for all *non-dummy* actions: C11 limits the *conditions* of an action to be only at *at start*, *over all* or *at end*. C12 limits the *effects* of an action to only happen *at start* or *at end*. In PDDL2.1 the structure of conditions/effects of all actions $\{a_j\}$ grounded from a particular operator are fixed. With this regard, C13 makes the conditions of all $\{a_j\}$ equal and C14 makes the effects of all $\{a_j\}$ equal. C15 makes the duration of all occurrences of the same action equal (if desired). Last but not least, C16 forces all actions to have at least one of its n -effects *at end*. Actions with only *at start* effects turn the value of the duration irrelevant besides they could exceed the plan makespan. (this last constraint is not specific of PDDL2.1 but produces more reasonable models for *durative actions*).

Table 3. Constraints to learn PDDL2.1-compliant action models.

ID	Constraint
C11.1	$(\text{req_start}(f, a) = \text{start}(a))$ OR $(\text{req_start}(f, a) = \text{end}(a))$
C11.2	$(\text{req_end}(f, a) = \text{start}(a))$ OR $(\text{req_end}(f, a) = \text{end}(a))$
C12	$(\text{time}(f, a) = \text{start}(a))$ OR $(\text{time}(f, a) = \text{end}(a))$
C13.1	$\forall f_i : (\forall a_j : \text{req_start}(f_i, a_j) = \text{start}(a_j))$ OR $(\forall a_j : \text{req_start}(f_i, a_j) = \text{end}(a_j))$
C13.2	$\forall f_i : (\forall a_j : \text{req_end}(f_i, a_j) = \text{start}(a_j))$ OR $(\forall a_j : \text{req_end}(f_i, a_j) = \text{end}(a_j))$
C14	$\forall f_i : (\forall a_j : \text{time}(f_i, a_j) = \text{start}(a_j))$ OR $(\forall a_j : \text{time}(f_i, a_j) = \text{end}(a_j))$
C15	$\forall a_i, a_j$ occurrences of the same action: $\text{dur}(a_i) = \text{dur}(a_j)$
C16	$\sum_{i=1}^n \text{time}(f_i, a) > n \times \text{start}(a)$

4.2.2 Mutex constraints

The set of mutexes that is given as input to a learning task \mathcal{L} allows to infer new information in form of *dynamic observations*: if two Boolean variables $\langle f_i, f_j \rangle$ are mutex they cannot hold simultaneously. This means that if we observe f_i , then we can infer $\neg f_j$ (despite $\neg f_j$ was not actually observed). This source of knowledge is specially relevant for the learning of *negative effects* when there is a lack of observations. Mutex information helps to fill this void by inferring the observation of negated variables, which forces later to satisfy the *causal links* of negative variables.

Given a $\langle f_i, f_j \rangle$ mutex, in our *temporal planning* model with *durative actions*, $\neg f_i$ does not necessarily implies f_j . See the effects (not (at ?t ?l1)) and (at ?t ?l2) of action *drive-truck* of Figure 1 that respectively happen *at start* and *at end*. If (at ?t ?l1) and (at ?t ?l2) are mutex (as defined in Figure 4), this means that the same truck cannot be in two locations simultaneously but that it is valid for the truck to be, for some time, at no location. These situations do not happen in STRIPS, where actions have instantaneous effects, so if $\langle f_i, f_j \rangle$ are mutex then f_i implies *not- f_j* and vice versa.

Mutex-constraints can be exploited in a pre-proces step for completing the input observations of a *one-shot learning task* \mathcal{L} . Furthermore, *dynamic observations* can be created to exploit mutex constraints at any generated intermediate state. This include states that where not observed but that are inferred by the CSP solutions. Given a mutex $\langle f_i, f_j \rangle$ it means that, immediately after a asserts f_i , we need to ensure the observation *not- f_j* . This is done while performing the CSP search, and if $\text{is_eff}(f_i, a)$ takes the value *true*, then the next observation is added: $\text{obs}(\text{not-}f_j, \text{time}(f_i, a) + \epsilon)$. The time of the observation cannot be just $\text{time}(f_i, a)$, as we first need to assert f_i and one ϵ later observe *not- f_j* . Adding the variables and constraints for this new observation is trivial for *Dynamic CSPs* (DCSPs). Otherwise, we need to statically define a new type of observation $\text{obs}(f_i, a, \text{not-}f_j)$, where a supports f_i which is mutex with f_j and, consequently, we will need to observe *not- f_j* . The difference w.r.t. an original obs is two-fold: i) the observation time is now initially unknown, and ii) the observation will be activated or not according to the following constraints:

if $(\text{is_eff}(f_i, a) = \text{true})$ **then** $(\text{start}(\text{obs}(f_i, a, \text{not-}f_j)) = \text{time}(f_i, a) + \epsilon)$ **AND**
 $(\text{is_cond}(\text{not-}f_j, \text{obs}(f_i, a, \text{not-}f_j)) = \text{true})$
else $\text{is_cond}(\text{not-}f_j, \text{obs}(f_i, a, \text{not-}f_j)) = \text{false}$

4.3 The CSP cost functions

The set of *conditions* of actions that are never deleted by any action are specially difficult to be learned with a pure satisfiability approach. For instance, the (link ?l1 ?l2) condition in the *drive-truck* action showed in the Figure 1. In general, this is an issue when learning action models in which *static predicates* appear in the action *conditions* [13].

This issue can be addressed extending the CP formulation to not only deal with the satisfaction of hard constraints but also to optimize a given *cost function* that defines the user preferences among the different possible action models that satisfy the given CSP preferring solutions that support the input observations in a way that is as *tight* as possible. To prefer this kind of *tight* support of the input observations we define the following two positive functions:

ϕ_1 *Initial causal-links*. This function counts the number of causal

links created to support the provided observations with plan actions. That is causal links $\langle b, f, a \rangle$ such that: (1) $\text{obs}(f, t)$ is an input observation and (2), action $b \neq \text{start}$, i.e., the supporter is not the *start* dummy action.

ϕ_2 *Side-effects*. This function counts the number of effects that are added by the actions and that do not build any causal link.

Our aim is to compute solutions to the CSP that minimizes both functions ϕ_1 and ϕ_2 . To achieve this we ask the CSP solve to *pareto optimize* functions ϕ_1 and ϕ_2 .

5 Planning, validation and learning with complete and incomplete action models

Here we show the flexibility of our CP formulation for addressing different task in the *temporal* (and *classical*) planning setting and with different amounts of input knowledge.

5.1 Complete and incomplete action models

The set of CSP variables $X3$, $X4$ and $X5$ from Table 1 (namely $\text{dur}(a)$, $\text{is_cond}(f, a)$ and $\text{is_eff}(f, a)$) represents the duration, the conditions and the effects of a given action a . If this information is known in advance for a given action a (e.g. because we are not learning from scratch but trying to complete a partially specified action model as introduced in 3) then these variables are set to the given known values. In this scenario the values of these variables can be propagated by the CSP reducing the branching factor of the solving process.

Further our formulation can be used straightforward as a similarity metric between durative actions. Given $\alpha(a)$, the *alphabet* (set of fluents) that can appear in the conditions and effects of a given durative action a , with parameters $\text{pars}(a)$ then the size of its space of possible action models is $\mathcal{D} \times 2^{5|\alpha(a)|}$ where \mathcal{D} is the number of different possible durations for a (in other words the domain of variable $X3$). Provided the *alphabet*, then the *conditions* and *effects* of a given durative action schema can be compactly coded by 5 bit-vectors, each of length $|\alpha(a)|$. A 0-bit in the vector represents that the corresponding *condition/effect* is not part of the schema while a 1-bit represents that is part of the schema. This also means that the *Hamming distance* can be used straightforward as a syntactic similarity metric between durative actions. For instance, we can use the *Hamming distance* to compare a learned action model with respect to a given reference model that serves as baseline. In this case, the number of wrong 1-bits in the learned schema provide us a measure of the *incorrectness* of the learned model (number of *conditions* and *effects* that should not be in the learned model) and the number of wrong 0-bits in the learned schema provide us a measure of the *incompleteness* of that model (number of *conditions* and *effects* that are missing in the learned model). A similar process was already defined for *strips* actions [1].

5.2 Integrating planning, validation and learning

Our CP formulation integrates the tasks of plan *synthesis*, plan *validation* and the learning of action models for the *temporal planning* setting. This connection lies on the fact that we can constrain the domain of the variables of our CP formulation to given known values. This feature is useful to leverage a priori knowledge of a given planning domain. For instance, because we have some available *prior knowledge* about the possible durations of a given action or because

we already know that a given action produces for sure certain effects or requires some conditions. In this case the value of the corresponding variables is a priori specified while the remaining variables are then regular variables whose value will be determined solving the CSP.

If all the variables that represent the conditions, effects and duration of the actions are a priori constrained to a single value (variables $X3$, $X4$ and $X5$) then solving the CSP is equivalent to solving a temporal planning task (that is synthesizing a plan that reaches a set of goals from certain initial state and with a given action model). Likewise, if all the variables that represent when the different actions appear in a solution plan (when the start times of actions happen, variables $X1$, $X2$ and $X3$) then solving the CSP is equivalent to validating a plan in a given temporal planning problem.

What is more, we can either synthesize (or validate) a plan despite some of the variables that representing the conditions, effects or duration of an action do not have a fixed value (its value is initially unknown). That is planning (and validating plans) when the action model is partially specified. Therefore, that the plan validation ability of our CP formulation is beyond the functionality of VAL (the standard plan validation tool [16]) since it can address plan validation of partial, or even empty, action models and with partially observed plan traces (VAL requires both a full plan and a full action model for plan validation).

The observations in \mathcal{O} can then be regarded as a sequence of ordered *landmarks* [15] for the planning problem $P(F, I, G, A?)$ since the fluents of the sets in \mathcal{O} must be achieved by any plan that solves $P_{\mathcal{O}}$ and in the same order as defined in the observation \mathcal{O} .

5.3 Classical planning

Last but not least this integration of the *planning*, *validation* and *learning* applies not only to the *temporal planning* setting but also to *classical planning*, the vanilla model of AI planning where actions are instantaneous [10]. In more detail, our CP formulation allows to transform the temporal planning model into a classical planning model by setting for each action $\text{dur}(a) = 0$ and constraining $\text{req_start}(a) = \text{req_end}(a) = \text{start}(a) = \text{time}(f, a)$

6 EVALUATION

[DE MOMENTO ESTO ESTA EN EL AIRE PORQUE NO SABEMOS COMO LO VAMOS A ABORDAR??]

The CP formulation has been implemented in Choco³, an open-source Java library for constraint programming that provides an object-oriented API to state the constraints to be satisfied. Choco uses a static model of variables and constraints, i.e. it is not a DCSP.

The empirical evaluation of a learning task can be addressed from two perspectives. From a pure syntactic perspective, learning can be considered as an automated design task to create a new model that is similar to a reference (or *ground truth*) model. Consequently, the success of learning is an accuracy measure of how similar these two models are, which usually counts the number of differences (in terms of incorrect durations or distribution of conditions/effects). Unfortunately, there is not a unique reference model when learning temporal models at real-world problems. Also, a pure syntax-based measure usually returns misleading and pessimistic results, as it may count as incorrect a different duration or a change in the distribution of conditions/effects that really represent equivalent reformulations of the

³ <http://www.choco-solver.org>

reference model. For instance, given the example of Figure 1, the condition learned (over all (link ?from ?to)) would be counted as a difference in action drive-truck, as it is at start in the reference model; but it is, semantically speaking, even more correct. Analogously, some durations may differ from the reference model but they should not be counted as incorrect. As seen in section ??, some learned durations cannot be granted, but the underlying model is still consistent. Therefore, performing a syntactic evaluation in learning is not always a good idea.

From a semantic perspective, learning can be considered as a classification task where we first learn a model from a training dataset, then tune the model on a validation test and, finally, asses the model on a test dataset. Our approach represents a one-shot learning task because we only use one plan sample to learn the model and no validation step is required. Therefore, the success of the learned model can be assessed by analyzing the success ratio of the learned model vs. all the unseen samples of a test dataset. In other words, we are interested in learning a model that fits as many samples of the test dataset as possible. This is the evaluation that we consider most valuable for learning, and define the success ratio as the percentage of samples of the test dataset that are consistent with the learned model. A higher ratio means that the learned model explains, or adequately fits, the observed constraints the test dataset imposes.

6.1 Learning from partially specified action models

We have run experiments on nine IPC planning domains. It is important to highlight that these domains are encoded in PDDL2.1, with the number of operators shown in Table 4, so we have included the constraints given in section 4.2.1. We first get the plans for these domains by using five planners (*LPG-Quality* [11], *LPG-Speed* [11], *TP* [17], *TFD* [6] and *TFLAP* [19]), where the planning time is limited to 100s. The actions and observations on each plan are automatically compiled into a CSP learning instance. Then, we run the one-shot learning task to get a temporal action model for each instance, where the learning time is limited to 100s on an Intel i5-6400 @ 2.70GHz with 8GB of RAM. In order to assess the quality of the learned model, we validate each model vs. the other models w.r.t. the *structure*, the *duration* and the *structure+duration*, as discussed in section 5. For instance, the *zenotravel* domain contains 78 instances, which means learning 78 models. Each model is validated by using the 77 remaining models, thus producing $78 \times 77 = 6006$ validations per struct, dur and struct+dur each. The value for each cell is the average success ratio. In *zenotravel*, the struct value means that the distribution of conditions/effects learned by using only one plan sample is consistent with all the samples used as dataset (100% of the 6006 validations), which is the perfect result, as also happens in *floortile* and *sokoban* domains. The dur value means the durations learned explain 68.83% of the dataset. This value is usually lower because any learned duration that leads to inconsistency in a sample counts as a failure. The struct+dur value means that the learned model explains entirely 35.76% of the samples. This value is always the lowest because a subtle structure or duration that leads to inconsistency in a sample counts as a failure. As seen in Table 4, the results are specially good, taking into consideration that we use only one sample to learn the temporal action model. These results depend on the domain size (number of operators, which need to be grounded), the relationships (causal links, threats and interferences) among the actions, and the size and quality of the plans.

We have observed that some planners return plans with unnecessary actions, which has a negative impact for learning precise du-

Table 4. Number of operators to learn. Instances used for validation. Average success ratio of the one-shot learned model vs. the test dataset in different IPC planning domains.

	ops	ins	struct	dur	struct+dur
<i>zenotravel</i>	5	78	100%	68.83%	35.76%
<i>driverlog</i>	6	73	97.60%	44.86%	21.04%
<i>depots</i>	5	64	55.41%	76.22%	23.19%
<i>rovers</i>	9	84	78.84%	5.35%	0.17%
<i>satellite</i>	5	84	80.74%	57.13%	40.53%
<i>storage</i>	5	69	58.08%	70.10%	38.36%
<i>floortile</i>	7	17	100%	80.88%	48.90%
<i>parking</i>	4	49	86.69%	81.38%	54.89%
<i>sokoban</i>	3	51	100%	87.25%	79.96%

rations. The worst result is returned in the *rovers* domain, which models a group of planetary rovers to explore the planet they are on. Since there are many parallel actions for taking pictures/samples and navigation of multiple rovers, learning the duration and the structure+duration is particularly complex in this domain.

6.2 Learning from scratch

7 CONCLUSIONS

REFERENCES

- [1] Diego Aineto, Sergio Jiménez, Eva Onaindia, and Miquel Ramírez, ‘Model recognition as planning’, in *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 29, pp. 13–21, (2019).
- [2] Eyal Amir and Allen Chang, ‘Learning partially observable deterministic action models’, *Journal of Artificial Intelligence Research*, **33**, 349–402, (2008).
- [3] S. N. Cresswell, T.L. McCluskey, and M.M West, ‘Acquiring planning domain models using LOCM’, *The Knowledge Engineering Review*, **28**(2), 195–213, (2013).
- [4] William Cushing, Subbarao Kambhampati, Daniel S Weld, et al., ‘When is temporal planning really temporal?’, in *Proceedings of the 20th international joint conference on Artificial intelligence*, pp. 1852–1859. Morgan Kaufmann Publishers Inc., (2007).
- [5] S. Edelkamp and J. Hoffmann, ‘PDDL2.2: the language for the classical part of IPC-4’, in *Proc. of the International Conference on Automated Planning and Scheduling (ICAPS-04) – International Planning Competition*, pp. 2–6, (2004).
- [6] Patrick Eyerich, Robert Mattmüller, and Gabriele Röger, ‘Using the context-enhanced additive heuristic for temporal and numeric planning’, in *Nineteenth International Conference on Automated Planning and Scheduling*, (2009).
- [7] Maria Fox and Derek Long, ‘PDDL2.1: An extension to PDDL for expressing temporal planning domains’, *Journal of artificial intelligence research*, **20**, 61–124, (2003).
- [8] Daniel Furelos Blanco, Antonio Bucciarone, and Anders Jonsson, ‘Carpool: Collective adaptation using concurrent planning’, in *AAMAS 2018. 17th International Conference on Autonomous Agents and Multi-agent Systems; 2018 Jul 10-15; Stockholm, Sweden.[Richland]: IFAA-MAS; 2018. International Foundation for Autonomous Agents and Multiagent Systems (IFAAMAS)*, (2018).
- [9] Antonio Garrido, Marlene Arangu, and Eva Onaindia, ‘A constraint programming formulation for planning: from plan scheduling to plan generation’, *Journal of Scheduling*, **12**(3), 227–256, (2009).
- [10] Hector Geffner and Blai Bonet, ‘A concise introduction to models and methods for automated planning’, *Synthesis Lectures on Artificial Intelligence and Machine Learning*, **8**(1), 1–141, (2013).
- [11] Alfonso Gerevini, Alessandro Saetti, and Ivan Serina, ‘Planning through stochastic local search and temporal action graphs in lpg’, *Journal of Artificial Intelligence Research*, **20**, 239–290, (2003).
- [12] Malik Ghallab, Dana Nau, and Paolo Traverso, *Automated Planning: theory and practice*, Elsevier, 2004.
- [13] Peter Gregory and Stephen Cresswell, ‘Domain model acquisition in the presence of static relations in the lop system’, in *Twenty-Fifth International Conference on Automated Planning and Scheduling*, (2015).
- [14] J. Hoffmann and S. Edelkamp, ‘The deterministic part of IPC-4: an overview’, *Journal of Artificial Intelligence Research*, **24**, 519–579, (2005).

- [15] Jörg Hoffmann, Julie Porteous, and Laura Sebastia, 'Ordered landmarks in planning', *Journal of Artificial Intelligence Research*, **22**, 215–278, (2004).
- [16] Richard Howey, Derek Long, and Maria Fox, '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*, pp. 294–301. IEEE, (2004).
- [17] Sergio Jiménez, Anders Jonsson, and Héctor Palacios, 'Temporal planning with required concurrency using classical planning', in *Proceedings of the 25th International Conference on Automated Planning and Scheduling (ICAPS)*, (2015).
- [18] Subbarao Kambhampati, 'Model-lite planning for the web age masses: The challenges of planning with incomplete and evolving domain models', in *Proceedings of the National Conference on Artificial Intelligence (AAAI-07)*, volume 22(2), pp. 1601–1604, (2007).
- [19] Eliseo Marzal, Laura Sebastia, and Eva Onaindia, 'Temporal landmark graphs for solving overconstrained planning problems', *Knowledge-Based Systems*, **106**, 14–25, (2016).
- [20] Kira Mourão, Luke S. Zettlemoyer, Ronald P. A. Petrick, and Mark Steedman, 'Learning STRIPS operators from noisy and incomplete observations', in *Conference on Uncertainty in Artificial Intelligence, UAI-12*, pp. 614–623, (2012).
- [21] Jussi Rintanen, 'Discretization of temporal models with application to planning with smt', in *Twenty-Ninth AAAI Conference on Artificial Intelligence*, (2015).
- [22] Vincent Vidal and Héctor Geffner, 'Branching and pruning: An optimal temporal poel planner based on constraint programming', *Artificial Intelligence*, **170**(3), 298–335, (2006).
- [23] Qiang Yang, Kangheng Wu, and Yunfei Jiang, 'Learning action models from plan examples using weighted MAX-SAT', *Artificial Intelligence*, **171**(2-3), 107–143, (2007).
- [24] Hankz Hankui Zhuo and Subbarao Kambhampati, 'Action-model acquisition from noisy plan traces', in *International Joint Conference on Artificial Intelligence, IJCAI-13*, pp. 2444–2450, (2013).
- [25] Hankz Hankui Zhuo and Subbarao Kambhampati, 'Model-lite planning: Case-based vs. model-based approaches', *Artificial Intelligence*, **246**, 1–21, (2017).
- [26] Hankz Hankui Zhuo, Tuan Nguyen, and Subbarao Kambhampati, 'Refining incomplete planning domain models through plan traces', in *Twenty-Third International Joint Conference on Artificial Intelligence*, (2013).
- [27] Hankz Hankui Zhuo, Tuan Anh Nguyen, and Subbarao Kambhampati, 'Refining incomplete planning domain models through plan traces', in *International Joint Conference on Artificial Intelligence, IJCAI-13*, pp. 2451–2458, (2013).
- [28] Hankz Hankui Zhuo, Qiang Yang, Derek Hao Hu, and Lei Li, 'Learning complex action models with quantifiers and logical implications', *Artificial Intelligence*, **174**(18), 1540–1569, (2010).