

One-Shot Learning of Concurrent Action Models

Antonio Garrido and Sergio Jiménez

Abstract. The paper presents a constraint programming (CP) formulation for the learning of models of *temporal planning* actions that can be executed in parallel and overlap. With the aim of understanding better the relations between the learning of action models, plan synthesis and plan validation, the paper focuses on the extreme learning scenario where just a single partial observation of the execution of a plan is available (i.e. one-shot). Like in existing CP formulations for temporal planning, our approach models time-stamps for actions, *causal link* relationships, *threats* and effect *interferences*. Further, we show that it can accommodate a different range of expressiveness, subsuming the PDDL2.1 temporal semantics, and that is solver-independent (i.e. off-the-shelf CSP solvers can be used for resolution).

1 Introduction

Temporal planning is an expressive planning model that relaxes the assumption of instantaneous actions of *classical planning* [8]. Actions in temporal planning are called *durative* (i.e. each action has an associated duration) and conditions/effects may hold/happen at different times [6]. This means that actions in the temporal planning model can be executed in parallel and overlap in several ways [4] and that valid solutions for temporal planning instances indicate the precise time-stamp when actions start and end [12].

Despite the potential of state-of-the-art planners, its applicability to the real world is still somewhat limited because of the difficulty of specifying correct and complete planning models [14]. The more expressive the planning model is, the more evident becomes this *knowledge acquisition bottleneck* that jeopardizes the usability of AI planning technology. This has led to a growing interest in the planning community for the learning of action models. Most approaches for learning planning action models are purely inductive and require large datasets of observations, e.g. thousands of plan observations to compute statistically significant models that minimizes some error metric over the observations [19, 17, 20, 15]. When model learning is considered an optimization task over a set of observations, it however does not guarantee that the learned action model may fail to explain a given input observation or that the states induced by the execution of plans generated with a learned model are incorrect.

This paper follows a different approach and analyzes the application of *Constraint Programming* (CP) for the *one-shot learning* of temporal action models, that is, for the extreme scenario where action models are learned from a single observation of the execution of a temporal plan. The contributions of this work are two-fold:

1. This is the first approach for learning action models for temporal planning where plan observations can refer to the execution of parallel and overlapping actions. This feature makes our approach appealing for learning in multi-agent environments. Learning classical action models from sequential plans is a well-studied
2. Our CP formulation relates the *learning* of planning action models with the *synthesis* and the *validation* of plans. Further, we show that the plan validation capacity of our CP formulation is beyond the functionality of VAL (the standard plan validation tool [12]) since it can address *plan validation* of partial (or even an empty) action models and with partially observed plan traces (VAL requires a full plan and a full action model for plan validation). Last but not least, our CP formulation allows that off-the-shelf CSP solvers can be used for the learning, planning and validation tasks.

As a motivating example, let us assume an observation of the executions of actions in a *logistics scenario*. Learning the temporal planning model will allow us: i) to better understand the insights of the logistics in terms of what is possible (or not) and why, because the model is consistent with the observed data; ii) to suggest changes that can improve the model originally created by a human, e.g. redistributing the actions' conditions, provided they still explain the observations; and iii) to automatically elaborate similar models for similar scenarios, such as public transit for commuters, tourists or people in general in metropolitan areas —*a.k.a.* smart urban mobility.

2 Background

This section formalizes the *temporal planning* model that we follow in this work, the hypothesis space for our learning target and the sampling space.

2.1 Temporal Planning

We assume that states are factored into a set F of Boolean variables. A state s is a time-stamped full assignment of values to these variables. The initial state I , is a fully observed state (i.e. $|I| = |F|$) time-stamped with $t = 0$. Last but not least let $G \subseteq F$ be a conjunction of goal conditions over F that defines the set of goal states.

A *temporal planning problem* is a tuple $\langle F, I, G, A \rangle$ where A represents the set of *durative actions*. Durative actions have an associated duration and may have conditions/effects at different times, such as conditions that must hold some time before the action starts,

effects that happen just when the action starts, in the middle of the action or some time after the action finishes [7].

In this work we assume that actions are grounded from *action schemes* (also known as *operators*) to compactly represent *temporal planning problems*. PDDL2.1 is a popular language for representing *temporal planning problems* and it is the input language for the temporal track of the International Planning Competition [6, 10]. PDDL2.1 somewhat restricts temporal expressiveness since it defines a *durative action* a with the following elements:

1. $\text{dur}(a)$, a positive value indicating the *duration* of the action.
2. $\text{cond}_s(a), \text{cond}_o(a), \text{cond}_e(a) \subseteq F$ representing the 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. In the simplest case, $\text{cond}_s(a) \cup \text{cond}_o(a) \cup \text{cond}_e(a) = \text{pre}(a)$.¹
3. $\text{eff}_s(a)$ and $\text{eff}_e(a)$ represent the action *effects*. In PDDL2.1 effects can happen *at start* or *at end* of a , respectively, and can still be positive or negative. Again, in the simplest case $\text{eff}_s(a) \cup \text{eff}_e(a) = \text{eff}(a)$.

The semantics of a PDDL2.1 durative action a can be defined in terms of two discrete events, $\text{start}(a)$ and $\text{end}(a) = \text{start}(a) + \text{dur}(a)$. This means that if action a starts on state s with time-stamp $\text{start}(a)$, $\text{cond}_s(a)$ 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 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.

A temporal plan is a set of pairs $\langle (a_1, t_1), (a_2, t_2) \dots (a_n, t_n) \rangle$. Each (a_i, t_i) pair contains a durative action a_i and a 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 a state $s_0 = I$, and $G \subseteq s_{\text{end}}$, being s_{end} the last state induced by the execution of the plan (therefore sequential plans can be expressed as temporal plans but not the opposite). We say that a temporal plan is a *solution* to a given temporal planning problem when the plan execution, starting from the corresponding initial state, eventually reaches a state s_{end} that meets the goal conditions.

2.2 The space of *durative* action models

In this work the learning target are FOL schemes of *durative* actions. Fig. 1 shows an example of two schemes for PDDL2.1 durative actions taken from the *driverlog* domain. The schema `board-truck` has a fixed duration while the duration of `drive-truck` depends on the driving time associated to the two given locations.

We assume that the set F of Boolean state variables is defined by the instantiation of a given set of predicates Ψ . We denote as $\mathcal{I}_{\xi, \Psi}$ the *vocabulary* (set of symbols) that can appear in the *conditions* and *effects* of a given *durative* action schema ξ . This set is defined as the FOL interpretations of the set of predicates Ψ , over the action parameters $\text{pars}(\xi)$.

For a *durative* action schema ξ , the size of its space of possible action models is then $D \times 2^{5 \times |\mathcal{I}_{\xi, \Psi}|}$ where D is the number of different possible durations for any action shaped by the ξ schema. This space

```
(: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 for PDDL2.1 durative actions.

is significantly larger than for learning STRIPS actions [19] where this number is just $2^{2 \times |\mathcal{I}_{\xi, \Psi}|}$ because negative effects must also be preconditions of the same action and cannot be positive effects of that action.

The conditions and effects of the schema of a *durative* action can be coded by 5 bit-vectors, each of length $|\mathcal{I}_{\xi, \Psi}|$. This means that the *Hamming distance* can be used straightforward as a similarity metric for *durative* schemes. For instance, to compare a learned action model with respect to a given reference model that serves as baseline. Note also that the Machine Learning metrics *precision* and *recall* can be also straightforward computed over the 5 bit-vectors that define the conditions and effect of the action schema of a *durative* action.

2.3 The sampling space

This paper studies the extreme scenario where just a single learning example is available (i.e. one-shot). In this work the single example is the noiseless but partial observation of the execution of a temporal plan starting from a given initial state.

This means that if the value of a fluent in F is observed then, that is the actual value of the fluent but that not the value of all the fluents can be observed at any time. Just a subset of them can be observed because, for instance, there are associated sensors reporting their value.

PDDL2.2 is an extension of PDDL2.1 that includes the notion of *Timed Initial Literal* [11] ($\text{til}(f, t)$), as a way of defining a fact $f \in F$ that becomes true at a certain time t , independently of the actions in the plan. Traditionally TILs are useful to define exogenous happenings; for instance, a time window when a warehouse is open in a logistics scenario, $\text{til}(\text{open}, 8)$ and $\text{til}(\neg \text{open}, 20)$. In this work we show that are also suitable to specify observations. The only differences is that as goals, observations represent conditions that must be satisfied in the execution of the plan at a particular time.

Figure 2. PDDL2.2 TILs representing the observation of a plan execution.

¹ Note that in classical planning, $\text{pre}(a) = \{p, \text{not} - p\}$ is contradictory. In temporal planning, $\text{cond}_s(a) = \{p\}$ and $\text{cond}_e(a) = \{\text{not} - p\}$ is a possible situation, though unusual

3 Learning Concurrent Action Models

This section formalizes the learning task we address in this paper and presents our CP formulation for addressing it with off-the-shelf CSP solvers.

3.1 One-shot learning of temporal action models

We define the task of the *one-shot learning of temporal action models* as the following tuple $\langle F, I, G, A?, O, C \rangle$, where:

- $\langle F, I, G, A? \rangle$ is a *temporal planning problem*, in which actions in $A?$ are partially specified (i.e., the exact conditions and effects and/or the duration of actions are unknown). In the worst case, we only know the vocabulary of the symbols that can appear in the *conditions* and/or *effects* of the actions. If prior knowledge is available we can use it to bound this vocabulary for a given action schema.
- O is the *observation* of a plan execution. At least it contains a full observation of the initial state (time-stamped with $t = 0$) and a time stamped final observation that equals the goals of the *temporal planning problem* $\langle F, I, G, A? \rangle$. In addition it can also contain time-stamped observations of the traversed intermediate states and about the time when actions started/ended their execution.
- C is a set of *state-constraints* that reflects domain-specific expert knowledge. These constraints allow us to complete the input observation and to prune the space of possible *durative* action models. Figure 3 show an example of a set of state-constraints for the *driverlog* domain.

$$\begin{aligned}
&\forall x_1, y_1, y_2 \neg at(x_1, y_1) \vee \neg at(x_1, y_2), \neq (y_1, y_2). \\
&\forall x_1, y_1, y_2 \neg in(x_1, y_1) \vee \neg in(x_1, y_2), \neq (y_1, y_2). \\
&\forall x_1, y_1, y_2 \neg driving(x_1, y_1) \vee \neg driving(x_1, y_2), \neq (y_1, y_2). \\
&\forall x_1, y_1, y_2 \neg driving(y_1, x_1) \vee \neg driving(y_2, x_1), \neq (y_1, y_2). \\
&\forall x_1, y_1, y_2 \neg at(x_1, y_1) \vee \neg driving(x_1, y_2). \\
&\forall x_1, y_1, y_2 \neg in(x_1, y_1) \vee \neg driving(x_1, y_2). \\
&\forall x_1, y_1, y_2 \neg at(x_1, y_1) \vee \neg in(x_1, y_2). \\
&\forall x_1, y_1 \neg empty(x_1) \vee \neg driving(y_1, x_1).
\end{aligned}$$

Figure 3. Examples of state-constraints for the *driverlog* domain.

A *solution* to the *one-shot learning of temporal action models* is a fully specified model of temporal actions \mathcal{A} such that: (1), the conditions, effects and duration of the actions is completely specified and (2), it is *consistent* with the given inputs $\langle F, I, G, A?, O, C \rangle$. In other words, for each action $a \in A?$, we have its equivalent version in \mathcal{A} where we have learned $dur(a)$, $cond_s(a)$, $cond_o(a)$, $cond_e(a)$, $eff_s(a)$ and $eff_e(a)$ and such that we can produce the actions in \mathcal{A} a plan whose execution starts in i , can produce the observations in O , and that reaches a final state that satisfies G .

3.2 Learning action models with constraint satisfaction

Our approach is to create a CSP whose solution induces an action model that solves a given *one-shot learning of temporal action models* task. Our CP formulation is inspired by previous work on *temporal planning as CP* [18, 7] and it is solver-independent. This means that any off-the-shelf CSP solver that supports the expressiveness of our CP formulation can be used.

3.2.1 The CSP Variables

For simplicity, we model time in \mathbb{Z}^+ and bound all maximum times to the makespan of the observed plan execution. If the observation of the plan makespan it is not available we consider a long enough domain for durations.

For each action a in $A?$, we create the seven kinds of variables for the CSP specified in Table 1. The value of the first three kind of CSP variables $start(a)$, $dur(a)$ and $end(a)$ and be given by the observation O or derived from the expression $end(a) = start(a) + dur(a)$. The remaining CSP variables model the interval when conditions must hold, the time when the effects happen and the causal links.

Variable	Domain	Description
$start(a)$	<i>Known/derived value</i>	start time for a
$end(a)$	<i>Known/derived value</i>	end time for a
$dur(a)$	$[1..max(a)]$	duration of a where $max(a) = makespan - start(a)$
$req_start(p, a)$, $req_end(p, a)$	$[0..makespan]$	interval $[req_start(p, a)..req_end(p, a)]$ during which action a requires p
$time(p, a)$	$[0..makespan]$	time when effect p of a happens
$sup(p, a)$	The action space	there is a causal link $\langle b_i, p, a \rangle$

Table 1. The CSP variables and their domains.

Additionally we create two dummy actions: *init* and *goal*. First, *init* represents the initial state I ($start(init) = 0$ and $dur(init) = 0$). *init* has no variables sup , req_start and req_end because it has no conditions. *init* has as many $time(p_i, init) = 0$ as p_i in I . Second, *goal* represents G ($start(goal) = makespan$ and $dur(goal) = 0$). *goal* has as many $sup(p_i, goal)$ and $req_start(p_i, goal) = req_end(p_i, goal) = makespan$ as p_i in G . *goal* has no variables time as it has no effects.

This formulation allows also to model TILs like any other regular actions. A $til(f, t)$ can be seen as a dummy action ($start(til(f, t)) = t$ and $dur(til(f, t)) = 0$) with no conditions and the single effect f that happens at time t ($time(f, til(f, t)) = t$). A *til* is analogous to *init*, as they both represent information that is given at a particular time, but externally to the execution of the plan.

3.2.2 The CSP Constraints

Table 2 shows the constraints defined among the CSP variables of Table 1. The first three constraints are explicit enough.

The fourth constraint models *causal links*. Note that in a causal link $\langle b, p, a \rangle$, $time(p, b) < req_start(p, a)$ and not \leq . This is because temporal planning assumes an $\epsilon > 0$ as a small tolerance between the time when a given effect p is supported and when it is required [6]. Since we model time in \mathbb{Z}^+ , $\epsilon = 1$ and \leq becomes $<$. The fifth constraint avoids any threat via promotion or demotion [10].

The sixth constraint models the fact that when the same action requires and deletes p then the effect cannot happen before the condition. Note the \geq inequality here; this is possible because if one condition and one effect of a happen at the same time, the underlying semantics in planning considers the condition is checked instantly before the effect [6].

The seventh constraint deals with the fact that two (possibly equal) actions have contradictory effects. It is important to note that these constraints involve any type of action, including *init* and *goal*.

In addition we can add the input *state-constraint* (if available), like the ones pictured in Figure 3, to our CSP formulation to prune the space of possible action models.

Constraint	Description
$\text{end}(a) = \text{start}(a) + \text{dur}(a)$	End time of a .
$\text{end}(a) \leq \text{start}(\text{goal})$	Always goal is the last action of the plan.
$\text{req_start}(p, a) \leq \text{req_end}(p, a)$	$[\text{req_start}(p, a) \dots \text{req_end}(p, a)]$ must be a valid interval.
if $\text{sup}(p, a) = b$ then $\text{time}(p, b) < \text{req_start}(p, a)$	Modeling causal link $\langle b, p, a \rangle$: the time when b supports p must be before a requires p .
$\forall c \neq a$ that deletes p at time t : if $\text{sup}(p, a) = b$ then $t < \text{time}(p, b)$ OR $t > \text{req_end}(p, a)$	Solving threat of c to causal link $\langle b, p, a \rangle$ being $c \neq a$ by promotion or demotion.
if a requires and deletes p : $\text{time}(\text{not} - p, a) \geq \text{req_end}(p, a)$	When a requires and deletes p , the effect cannot happen before the condition.
$\forall a_i, a_j \mid a_i$ supports p and a_j deletes p : $\text{time}(p, a_i) \neq \text{time}(\text{not} - p, a_j)$	Solving effect interference (p and $\text{not} - p$): they cannot happen at the same time.

Table 2. The CSP constraints.

3.3 Specific constraints for the PDDL2.1 action model

This temporal model formulation is more expressive than PDDL2.1 (see more details in section 3.3), and allows conditions and effects to be at any time, even outside the execution of the action. For instance, let us imagine a condition p that only needs to be maintained for 5 time units before an action a starts (e.g. warming-up a motor before driving): the expression $\text{req_end}(p, a) = \text{start}(a)$; $\text{req_end}(p, a) = \text{req_start}(p, a) + 5$ is possible in our formulation. Additionally, we can represent an effect p that happens in the middle of action a : $\text{time}(p, a) = \text{start}(a) + (\text{dur}(a)/2)$ is also possible.

We present now our formulation for PDDL2.1. As section 2.1 explains, PDDL2.1 restricts the expressiveness of temporal planning in terms of conditions, effects, durations and structure of the actions. Hence, our temporal formulation subsumes and is significantly richer than PDDL2.1; but adding constraints to make it fully PDDL2.1-compliant is straightforward.

First, adding $((\text{req_start}(p, a) = \text{start}(a)) \text{ OR } (\text{req_start}(p, a) = \text{end}(a))) \text{ AND } ((\text{req_end}(p, a) = \text{start}(a)) \text{ OR } (\text{req_end}(p, a) = \text{end}(a)))$ limits condition p to be *at start*, *over all* or *at end*, i.e. p is in $\text{cond}_s(a)$, $\text{cond}_o(a)$ or $\text{cond}_e(a)$, respectively. Further, if a condition is never deleted in a plan, it can be considered as an invariant condition for such a plan. In other words, it represents static information. This type of condition is commonly used in planning to ease the grounding process from the operators; e.g. to model that there is a link between two locations and, consequently, a driving is possible, or modeling a petrol station that allows a refuel action in a given location, etc. Therefore, the constraint to be added for an invariant condition p is simply: $((\text{req_start}(p, a) = \text{start}(a)) \text{ AND } (\text{req_end}(p, a) = \text{end}(a)))$, i.e. $p \in \text{cond}_o(a)$. Surprisingly, invariant conditions are modeled differently depending on the human modeler. See, for instance, $(\text{link} \text{ ?from } ?\text{to})$ of Fig. 1, which is modeled as an *at start* condition despite: i) the link should be necessary all over the driving, and ii) no action in this domain can be planned to delete that link. This also happens in the *transport* domain of the IPC, where a refuel action requires to have a petrol station in a location only *at start*, rather than *over all* which makes more sense. This shows that modeling a planning domain is not easy and it highly depends on the human’s decision. On the contrary, our formulation checks the invariant conditions and deals with them always in the same coherent way.

Second, $((\text{time}(p, a) = \text{start}(a)) \text{ OR } (\text{time}(p, a) = \text{end}(a)))$ makes an effect p happen only *at start* or *at end* of action a , i.e. p

is in $\text{eff}_s(a)$ or $\text{eff}_e(a)$. Also, if all effects happen *at start* the duration of the action would be irrelevant and could exceed the plan makespan. To avoid this, for any action a , at least one of its effects should happen *at end*: $\sum_{i=1}^{n=|\text{eff}(a)|} \text{time}(p_i, a) > n \cdot \text{start}(a)$, which guarantees $\text{eff}_e(a)$ is not empty.

Third, durations in PDDL2.1 can be defined in two different ways. On the one hand, durations can be equal for all grounded actions of the same operator. For instance, any instantiation of `board-truck` of Fig. 1 will last 2 time units no matter its parameters. Although this may seem a bit odd, it is not an uncommon practice to simplify the model. The constraint to model this is: $\forall a_i, a_j$ being instances of the same operator: $\text{dur}(a_i) = \text{dur}(a_j)$. On the other hand, although different instantiations of `drive-truck` will last different depending on the locations, different occurrences of the same instantiated action will last equal. In a PDDL2.1 temporal plan, multiple occurrences of `drive-truck` (`truck1, loc1, loc2, driver1`) will have the same duration no matter when they start. Intuitively, they are different occurrences of the same action, but in the real-world the durations would differ from driving at night or in peak times. Since PDDL2.1 makes no distinction among different occurrences, the constraint to add is: $\forall a_i, a_j$ being occurrences of the same durative action: $\text{dur}(a_i) = \text{dur}(a_j)$. Obviously, this second constraint is subsumed by the first one in the general case where all instances of the same operator have the same duration.

Fourth, the structure of conditions and effects for all grounded actions of the same operator is constant in PDDL2.1. This means that if $(\text{empty} \text{ ?t})$ is an *at start* condition of `board-truck`, it will be *at start* in any of its grounded actions. Let $\{p_i\}$ be the conditions of an operator and $\{a_j\}$ be the instances of a particular operator. The following constraints are necessary to guarantee a constant structure:

$$\forall p_i : (\forall a_j : \text{req_start}(p_i, a_j) = \text{start}(a_j)) \text{ OR } (\forall a_j : \text{req_start}(p_i, a_j) = \text{end}(a_j))$$

$$\forall p_i : (\forall a_j : \text{req_end}(p_i, a_j) = \text{start}(a_j)) \text{ OR } (\forall a_j : \text{req_end}(p_i, a_j) = \text{end}(a_j))$$

And analogously for all effects $\{p_i\}$ and the instances $\{a_j\}$ of an operator:

$$\forall p_i : (\forall a_j : \text{time}(p_i, a_j) = \text{start}(a_j)) \text{ OR } (\forall a_j : \text{time}(p_i, a_j) = \text{end}(a_j))$$

As a conclusion, in our formulation each action of $A?$ is modeled separately so it does not need to share the same structure or duration of other actions. Moreover, the time-stamps for conditions/effects can be arbitrarily placed inside or outside the execution of the action, which allows for a flexible and expressive temporal model. But, when necessary, we can simply include additional constraints to restrict the expressiveness of the model, such as the ones provided by PDDL2.1.

3.4 Example

We now show a fragment of the formulation for the example depicted in Fig. ?? . For simplicity, we only show the variables and constraints for action $a3$, but the formulation is analogous for all other actions.

The variables and domains are: $\text{start}(a3) = 7$; $\text{dur}(a3) \in [1..15]$; $\text{end}(a3) = \text{start}(a3) + \text{dur}(a3)$; $\text{sup}(r, a3) \in \{a1\}$; $\text{req_start}(r, a3), \text{req_end}(r, a3) \in [0..22]$; and $\text{time}(\text{not} - p, a3) \in [0..22]$. On the other hand, the constraints are: $\text{end}(a3) \leq \text{start}(\text{goal})$; $\text{req_start}(r, a3) \leq \text{req_end}(r, a3)$; if $\text{sup}(r, a3) = a1$ then $\text{time}(r, a1) < \text{req_start}(r, a3)$; if $\text{sup}(r, a3) = a1$ then $((\text{time}(\text{not} - r, a4) < \text{time}(r, a1)) \text{ OR } (\text{time}(\text{not} - r, a4) > \text{req_end}(r, a3)))$; $\text{time}(\text{not} - p, a3) \neq \text{time}(p, a1)$ and $\text{time}(\text{not} - p, a3) \neq \text{time}(p, a2)$.

There are many consistent solutions for this simple example, mainly because there is a huge range of possible durations that make the learned model consistent with the partially specified model $A?$. Fig. 4 shows six arbitrary models as solutions. What is important to note is that the structure, i.e. distribution of conditions/effects, is similar in all the learned models. Actually, the distribution of the effects is identical (except for q in model 2), and the distribution of conditions is very similar (e.g. q is always in cond_o and r in $a4$ is very often in cond_o). This shows that the one-shot learning returns not only consistent models but also similar, which is very positive. The durations are, however, more different: $\text{dur}(a1)$ ranges in these models from 7 to 19, whereas $\text{dur}(a2)$ ranges from 5 to 18. As explained in section ??, learning the precise duration from just one sample may not be always possible, which is the main limitation of the one-shot learning task. In fact, the specific constraint of PDDL2.1, with regard to having multiple occurrences of the same action having the same duration, can significantly help us to learn the actions' duration in a more precise way as the learned duration must be consistent with all those occurrences.

Action	dur	cond _s	cond _o	cond _e	eff _s	eff _e
Learned model 1						
a1	8				r	p
a2	18				q	p
a3	1		r			$\text{not} - p$
a4	1		q, r	p		$\text{not} - r$
Learned model 2						
a1	19				r	p
a2	5					p, q
a3	1			r		$\text{not} - p$
a4	1	r	p, q			$\text{not} - r$
Learned model 3						
a1	7				r	p
a2	18				q	p
a3	1			r		$\text{not} - p$
a4	1		p, q, r			$\text{not} - r$

Figure 4. Three learned models for the example of Fig. ??.

3.5 Implementation. Use of Heuristics for Resolution

Our CSP formulation is automatically compiled from a partially specified action model, as defined in a classical planning problem, and the observations from a plan execution. The formulation has been implemented in Choco², an open-source Java library for constraint programming that provides an object-oriented API to state the constraints to be satisfied.

Our formulation is solver-independent, which means we do not use heuristics that may require changes in the implementation of the CSP engine. Although this reduces the solver performance, we are interested in using it as a blackbox that can be easily changed with no modification in our formulation. However, we can easily encode standard static heuristics for variable and value selection that help improve efficiency by following the next ordering, which has shown very efficient in our experiments:

1. Effects (time). For negative effects, first the lower value and for positive effects, first the upper value. This gives priority to delete effects as $\text{eff}_s(a)$ and positive effects as $\text{eff}_e(a)$.

² <http://www.choco-solver.org>

2. Conditions (req_start and req_end). For req_start, first the lower value, whereas for req_end, first the upper value. This gives priority to $\text{cond}_o(a)$, trying to keep the conditions as long as possible.
3. Supporters (sup). First the lower value, thus preferring the supporter that starts earlier in the plan.
4. Duration (dur). First the lower value, thus applying the principle of the shortest actions that make the learned model consistent.

3.6 Using the CP Formulation for Plan Validation

The *one-shot learning of temporal action models* is related to plan *synthesis* and plan *validation*: In the case of plan synthesis the intermediate observations in O are empty while in the case of plan validation the actions in $A?$ are fully specified.

We explained that adding extra constraints allows us to restrict the temporal expressiveness of the learned model. We show here that we can also restrict the learned model by constraining the variables to known values, which is specially interesting when there is additional information on the temporal model that needs to be represented. For instance, based on past learned models, we may know the precise duration of an action a is 6, or we can figure out that an effect p always happens at end. Our CP formulation can include this by simply adding $\text{dur}(a) = 6$ and $\text{time}(p, a) = \text{end}(a)$, respectively, which is useful to enrich the partially specified actions in $A?$ of the learning task.

In particular, the possibility of adding those constraints is very appealing when used for validating whether a partial action model allows us to learn a consistent model, as we will see in section 4. Let us assume that the distribution of all (or just a few) conditions and/or effects is known and, in consequence, represented in the learning task. If a solution is found, then that structure of conditions/effects is consistent for the learned model. On the contrary, if no solution is found that structure is inconsistent and cannot be explained. Analogously, we can represent known values for the durations. If a solution is found, the durations are consistent, and inconsistent otherwise. Hence, we have three options for validating a partial model *w.r.t.*: i) a known structure with the distribution of conditions/effects; ii) a known set of durations; and iii) a known structure plus a known set of durations (i+ii). The first and second option allows for some flexibility in the learning task because some variables remain open. On the contrary, the third option checks whether a learned model can fit the given constraints, thus reproducing a strict plan validation task equivalent to [12].

4 Evaluation

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 reference model. For instance, given the example of Fig. 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, assess 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.

4.1 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 3, so we have included the constraints given in section 3.3. We first get the plans for these domains by using five planners (*LPG-Quality* [9], *LPG-Speed* [9], *TP* [13], *TFD* [5] and *TFLAP* [16]), 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 3.6. 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 3, 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 durations. 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

	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%

Table 3. 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.

and navigation of multiple rovers, learning the duration and the structure+duration is particularly complex in this domain.

4.2 Learning from scratch

5 Conclusions

We have presented a purely declarative CP formulation, which is independent of any CSP solver, to address the learning of temporal action models. Learning in planning is specially interesting to recognize past behavior in order to predict and anticipate actions to improve decisions. The main contribution is a simple formulation that is automatically derived from the actions and observations on each plan execution, without the necessity of specific hand-coded domain knowledge. It is also flexible to support a very expressive temporal planning model, though it can be easily modified to be PDDL2.1-compliant. Formal properties are inherited from the formulation itself and the CSP solver. The formulation is correct because the definition of constraints to solve causal links, threats and effect interferences are supported, which avoids contradictions. It is also complete because the solution needs to be consistent with all the imposed constraints, while a complete exploration of the domain of each variable returns all the possible learned models in the form of alternative consistent solutions.

Unlike other approaches that need to learn from datasets with many samples, we perform a one-shot learning. This reduces both the size of the required datasets and the computation time. The one-shot learned models are very good and explain a high number of samples in the datasets used for testing. Moreover, the same CP formulation is valid for learning and for validation, by simply adding constraints to the variables. This is an advantage, as the same formulation allows us to carry out different tasks: from entirely learning, partial learning/validation (structure and/or duration) to entirely plan validation. According to our experiments, learning the structure of the actions in a one-shot way leads to representative enough models, but learning the precise durations is more difficult, and even impossible, when many actions are executed in parallel.

Finally, our CP formulation can be represented and solved by Satisfiability Modulo Theories, which is part of our current work. As future work, we want to extend our formulation to learn meta-models, as combinations of many learned models, and a more complete action model. In the latter, rather than using a partially specified set of actions, we want to find out the conditions/effects together with their distribution. The underlying idea of finding an action model consistent with all the constraints will remain the same, but the model will need to be extended with additional decision variables and constraints. This will probably lead to the analysis of new heuristics for resolution.

REFERENCES

- [1] Eyal Amir and Allen Chang, 'Learning partially observable deterministic action models', *Journal of Artificial Intelligence Research*, **33**, 349–402, (2008).
- [2] Ankuj Arora, Humbert Fiorino, Damien Pellier, Marc Métivier, and Sylvie Pesty, 'A review of learning planning action models', *The Knowledge Engineering Review*, **33**, (2018).
- [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] 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).
- [6] 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).
- [7] 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).
- [8] 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).
- [9] 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).
- [10] Malik Ghallab, Dana Nau, and Paolo Traverso, *Automated Planning: theory and practice*, Elsevier, 2004.
- [11] J. Hoffmann and S. Edelkamp, 'The deterministic part of IPC-4: an overview', *Journal of Artificial Intelligence Research*, **24**, 519–579, (2005).
- [12] 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).
- [13] 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).
- [14] 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).
- [15] Jirí Kucera and Roman Barták, 'LOUGA: learning planning operators using genetic algorithms', in *Pacific Rim Knowledge Acquisition Workshop, PKAW-18*, pp. 124–138, (2018).
- [16] Eliseo Marzal, Laura Sebastia, and Eva Onaindia, 'Temporal landmark graphs for solving overconstrained planning problems', *Knowledge-Based Systems*, **106**, 14–25, (2016).
- [17] 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).
- [18] Vincent Vidal and Héctor Geffner, 'Branching and pruning: An optimal temporal pocl planner based on constraint programming', *Artificial Intelligence*, **170**(3), 298–335, (2006).
- [19] 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).
- [20] 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).
- [21] Hankz Hankui Zhuo and Subbarao Kambhampati, 'Model-lite planning: Case-based vs. model-based approaches', *Artificial Intelligence*, **246**, 1–21, (2017).
- [22] 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).
- [23] 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).