



# Computational models of planning

Hector Geffner\*

The selection of the action to do next is one of the central problems faced by autonomous agents. Natural and artificial systems address this problem in various ways: action responses can be hardwired, they can be learned, or they can be computed from a model of the situation, the actions, and the goals. Planning is the model-based approach to action selection and a fundamental ingredient of intelligent behavior in both humans and machines. Planning, however, is computationally hard as the consideration of all possible courses of action is not computationally feasible. The problem has been addressed by research in Artificial Intelligence that in recent years has uncovered simple but powerful computational principles that make planning feasible. The principles take the form of domain-independent methods for computing heuristics or appraisals that enable the effective generation of goal-directed behavior even over huge spaces. In this paper, we look at several planning models, at methods that have been shown to scale up to large problems, and at what these methods may suggest about the human mind. © 2013 John Wiley & Sons, Ltd.

## How to cite this article:

WIREs Cogn Sci 2013, 4:341–356. doi: 10.1002/wcs.1233

## INTRODUCTION

Planning is one of the key aspects of intelligent behavior, and one of the first to be studied in the field of Artificial Intelligence (AI). The first AI planner and one of the first AI programs was the General Problem Solver (GPS) developed by Newell and Simon in the late 1950s.<sup>1,2</sup> Since then, planning has remained a central topic in AI while changing in significant ways: on the one hand, it has become more rigorous, with a variety of planning models defined and studied; on the other, it has become more empirical, with planning algorithms evaluated experimentally and a strong emphasis placed on scalability: effective goal-oriented behavior that is not limited by the number of actions and variables in the problem.

Planning can be understood as representing one of the three main approaches for *selecting the action to do next*; a problem that is central in the design of autonomous systems. In the *hardwired approach*, action responses are hardwired by nature, or more commonly in AI systems, by a programmer. For a

robot moving in an office environment, for example, the program may say to back up when too close to a wall, to search for a door if the robot has to move to another room, and so on.<sup>3,4</sup> The problem with this approach, common to many other AI systems, is that since it is difficult to anticipate all possible situations and the best way to handle them, the resulting systems tend to be brittle. In the *learning-based approach*, on the other hand, action responses are not hardwired but are learned by trial and error as in reinforcement learning,<sup>5</sup> or by generalization from examples as in supervised learning.<sup>6,7</sup> Finally, in the *model-based approach*, action responses are computed from a model of the situation, the actions, the sensors, and goals.<sup>7–9</sup> The distinction that the philosopher Daniel Dennett makes between ‘Darwinian’, ‘Skinnerian’, and ‘Popperian’ creatures,<sup>10</sup> mirrors quite closely the distinction between hardwired agents, agents that learn, and agents that use models respectively. The approaches, however, are not incompatible, and indeed, agents that learn models rather than behavioral responses fit also in the latter class, as they must then use the models learned for selecting the action to do next.

Planning, usually associated with ‘thinking before acting’, is best conceived as the model-based

\*Correspondence to: hector.geffner@upf.edu

ICREA & Universitat Pompeu Fabra, Roc Boronat 138, Barcelona, Spain

The authors have declared no conflicts of interest for this article.

approach to action selection. It is a fundamental ingredient of intelligent behavior in both humans and machines, yet it is computationally hard. Indeed, even in the most simple planning model where the state of the world is fully known and actions have deterministic effects, the number of possible world states is exponential in the number of state variables, and the number of possible plans is exponential in the length of the plans.<sup>a</sup> The situation facing an artificial agent that plans even in the simplest setting, has a key element that is common with humans making plans: neither one can afford to consider all possible courses of action as this is not computationally feasible. This computational challenge has been used as evidence for contesting the possibility of general planning and reasoning abilities in humans or machines.<sup>12</sup> The complexity of planning, however, just implies that no planner can efficiently solve every problem from every domain, not that a planner cannot solve an infinite collection of problems from old and new domains, and hence be useful to an acting agent. This is indeed the way modern AI planners are empirically evaluated and ranked in AI planning competitions, where they are tried on domains that the planners' authors have never seen.<sup>13</sup> Thus, far from representing an insurmountable obstacle, the twin requirements of generality and scalability have been addressed head on in AI planning research, and this has led to simple but powerful computational principles that make domain-independent planning feasible. The principles take the form of heuristics or appraisals that are computed automatically for the problem at hand from suitable relaxations (simplifications) of the problem. The heuristics derived in this way tailor the general solver to the specific problem, and enable the generation of goal-directed behavior even over huge spaces. In this paper, we look at these methods and what they may reveal about the human mind, turning the requirements of generality and scalability from an impossibility into a critical source of insight.

From a cognitive perspective, one lesson to draw is that if one of the primary roles of feelings and emotions is to act as rough guides in the generation of behavior in complex physical and social worlds where the consideration of all possible courses of action is not computationally feasible, there is much to be gained by looking at suitable abstractions of the problem that present the same computational challenges in crisp form. More specifically, we will argue that the domain-independent methods for deriving heuristic functions that have been developed for making planning feasible provide a different angle from which to look at the computation and role of emotional appraisals in humans.<sup>14,b</sup>

The paper is organized as follows. We consider the most basic planning model first, the computational challenge that it presents, and the domain-independent methods and heuristics that have been developed for addressing this challenge. We then look at what these domain-independent heuristics suggest about the generation, nature, and role of unconscious appraisals, and move on to alternative planning methods and richer planning models. We conclude by discussing some open challenges in planning research and summarizing the main lessons.

## BASIC MODEL: CLASSICAL PLANNING

The most basic model in planning is concerned with the selection of actions for achieving goals when the initial situation is fully known and actions have deterministic effects. The model underlying this form of planning, called *classical planning*, can be described formally in terms of a state space featuring

- a finite and discrete set of states  $S$ ,
- a *known initial state*  $s_0 \in S$ ,
- a set  $S_G \subseteq S$  of goal states,
- a set of actions  $A(s) \subseteq A$  applicable in each state  $s \in S$ ,
- a *deterministic state transition function*  $f(a, s)$  for  $a \in A(s)$  and  $s \in S$ , so that  $f(a, s)$  denotes the state that results from doing action  $a$  in state  $s$ , and
- positive *action costs*  $c(a, s)$  that may depend on the action  $a$  and the state  $s$ .

A solution or *plan* for this model is a sequence of actions  $a_0, \dots, a_n$  that generates a state sequence  $s_0, s_1, \dots, s_{n+1}$  such that each action  $a_i$  is applicable in the state  $s_i$  and results in the state  $s_{i+1} = f(a_i, s_i)$ , the last of which is a goal state; that is  $s_{i+1} = f(a_i, s_i)$  for  $a_i \in A(s_i)$ ,  $i = 0, \dots, n$ , and  $s_{n+1} \in S_G$ . The cost of a plan is the sum of the action costs  $c(a_i, s_i)$ , and a plan is optimal if it has minimum cost. The cost of a problem is the cost of its optimal solutions. When action costs are all one, plan cost reduces to plan length, and the optimal plans are simply the shortest ones.

As an illustration, the problem of rearranging a set of block towers into a different set of towers can be naturally formulated as a model of this type where the state represents the possible block configurations. Similarly, delivering a set of goods to clients, solving a puzzle such as Rubik's Cube, or playing solitaire can be all expressed as classical planning problems.

The model underlying classical planning does not account for either uncertainty or sensing, and thus

gives rise to plans that represent open-loop controllers where observations play no role. Planning models that take these aspects into account give rise to different types of controllers and will be considered later.

## STATE VARIABLES AND FACTORED REPRESENTATIONS

The state models required for planning are not represented explicitly in general because they are often too large. The planning agent is assumed to have instead a compact and implicit representation of the state model given in terms of a set of state variables. The states are the possible combination of values over these variables, and the actions are defined in terms of pre and postconditions expressed over the state variables as well.

One of the most common languages for representing classical planning problems is also one of the oldest, Strips,<sup>17</sup> a planning language that can be traced back to the late 1960s. A planning problem in Strips is a tuple  $P = \langle F, O, I, G \rangle$  where

- $F$  represents a set of *Boolean variables*,
- $O$  represents a set of *actions*,
- $I$  represents the *initial situation*, and
- $G$  represents the *goal*.

Both the initial situation  $I$  and the goal  $G$  are expressed by a set of atoms over  $F$ . For a Boolean state variable  $p$  in  $F$ , there are two atoms  $p = \text{true}$  and  $p = \text{false}$ , abbreviated using logical notation as  $p$  and  $\neg p$ , where the symbol ‘ $\neg$ ’ stands for negation. A state in a Strips problem is a valuation of the state variables represented by the collection of atoms that the state makes true. The initial situation  $I$  stands for the state that makes the atoms in  $I$  true and all other atoms false, while the goal  $G$  stands for the collection of states that make all the atoms in  $G$  true.

In Strips, the actions  $o \in O$  are represented by three sets of atoms over  $F$  called the Add, Delete, and Precondition lists, denoted as  $\text{Add}(o)$ ,  $\text{Del}(o)$ ,  $\text{Pre}(o)$ . The first describes the atoms that the action  $o$  makes true, the second, the atoms that  $o$  makes false, and the third, the atoms that must be true in order for the action to be applicable. A Strips problem  $P = \langle F, O, I, G \rangle$  encodes implicitly, in compact form, the classical state model  $S(P)$  where

- the states  $s \in S$  are the possible *collections of atoms* over  $F$ ,
- the initial state  $s_0$  is  $I$ ,

- the goal states  $s$  are those for which  $G \subseteq s$ ,
- the actions  $a$  in  $A(s)$  are the ones in  $O$  with  $\text{Pre}(a) \subseteq s$ ,
- the state transition function is  $f(a, s) = (s \setminus \text{Del}(a)) \cup \text{Add}(a)$ , that is, the state  $s$  with the atoms in  $\text{Del}(a)$  deleted and the atoms in  $\text{Add}(a)$  added, and
- the action costs  $c(a)$  are equal to one by default.

Given that the Strips problem represents the state model  $S(P)$ , the *plans* for  $P$  are defined as the plans for  $S(P)$ ; namely, the action sequences that map the initial state  $s_0$  that corresponds to  $I$  into a goal state where the goals  $G$  are true. Since the states in  $S(P)$  are represented as collections of atoms from  $F$ , the number of states in  $S(P)$  is  $2^{|F|}$  where  $|F|$  is the number of Boolean variables  $F$  in  $P$ , usually called *fluents*.

The state representation that follows from a planning language such as Strips is domain-independent. Thus, while a specialized solver for the Blocks World may represent the state by a set of lists, each one representing a tower of blocks, in the state representation that follows from Strips there will be just atoms, and the same will be true of any other domain. As an illustration, a domain that involves three locations  $l_1, l_2$ , and  $l_3$ , and three tasks  $t_1, t_2$ , and  $t_3$ , where  $t_i$  can be performed only at location  $l_i$ , can be modeled with a set  $F$  of fluents  $\text{at}(l_i)$  and  $\text{done}(t_i)$ , and a set  $O$  of actions  $\text{go}(l_i, l_j)$  and  $\text{do}(t_i)$ ,  $i, j = 1, \dots, 3$ , with precondition, add, and delete lists

$$\text{Pre}(a) = \{\text{at}(l_i)\}, \text{Add}(a) = \{\text{at}(l_j)\}, \text{Del}(a) = \{\text{at}(l_i)\}$$

for  $a = \text{go}(l_i, l_j)$ , and

$$\text{Pre}(a) = \{\text{at}(l_i)\}, \text{Add}(a) = \{\text{done}(t_i)\}, \text{Del}(a) = \{\}$$

for  $a = \text{do}(t_i)$ . The problem of doing tasks  $t_1$  and  $t_2$  starting at location  $l_3$  can then be modeled by the tuple  $P = \langle F, I, O, G \rangle$  where

$$I = \{\text{at}(l_3)\} \text{ and } G = \{\text{done}(t_1), \text{done}(t_2)\}.$$

A solution to  $P$  is an applicable action sequence that maps the state  $s_0 = I$  into a state where the goals in  $G$  are all true. In this case one such plan is the sequence

$$\pi = \{\text{go}(l_3, l_1), \text{do}(t_1), \text{go}(l_1, l_2), \text{do}(t_2)\}.$$

The number of states in the problem is  $2^6$  as there are six Boolean variables. Still, it can be shown that many of these states are not reachable from the initial state. Indeed, the atoms  $\text{at}(l_i)$  for  $i = 1, 2, 3$  are



be expressed as a Strips problem  $P = \langle F, I, O, G \rangle$  with a set of atoms  $F$  given by  $on(x, y)$ ,  $ontable(x)$ , and  $clear(x)$ , where  $x$  and  $y$  range over the block labels  $A$ ,  $B$ , and  $C$ . The figure shows the graph associated to the problem, whose solution is a path connecting the node representing the initial situation with a node representing a goal situation.

The Blocks World is simple for people but until recently, not so easy for *domain-independent* planners that must accept *any* planning problem  $P = \langle F, I, O, G \rangle$  in their input, for *any domain*, and solve it automatically *without assuming additional knowledge*. Indeed, the number of states in a Blocks World problem with  $n$  blocks is exponential in  $n$ , as the states include all the  $n!$  possible towers of  $n$  blocks plus additional combinations of lower towers.

It may be argued that people solve a problem like Blocks World with *domain-specific control knowledge*; namely, knowledge about *how* problems in the domain are solved, rather than just knowledge about *what* the problem is about. Domain-specific control knowledge for Blocks, for example, may state that ‘bad placed’ blocks and blocks above them must be all moved, and that blocks should be moved onto other blocks only when the latter are ‘well placed’. Many arguments have been raised against the view of humans as ‘general problem solvers’, some casting the mind as a collection of specialized modules evolved to solve the specific problems of our Pleistocene ancestors.<sup>12</sup> Yet, the arguments against generality are speculative, and it is not clear what the modules are, how specialized they are, nor what are the specific problems that they evolved to solve. Also, if there are many modules, some mechanism would still have to determine which module to use and when, and this mechanism cannot be domain-specific. Approaches that deny general problem solving abilities may be actually pushing the complexity of the problem one level up, to the module or modules that must control the rest of the modules.

Actually, these questions involve computational and complexity issues that cannot be answered fully without knowing whether some form of ‘generality’ is computationally feasible, and what the price for this generality is. Indeed, why deny general problem solving abilities a priori if it is computationally feasible and cheap? Of course, because of the theoretical complexity of planning we cannot expect a general planner to efficiently solve every problem from every possible domain, yet a general planner will be useful enough for an agent if it can solve problems from many domains, provided that the domains themselves are not inherently that hard, and that the size of a problem, as measured by the number of actions and

variables, is not by itself an impediment to its solution. This is actually what modern AI planners manage to do. We will see below how this is accomplished and why this is relevant from a cognitive perspective.

## HEURISTICS

One way to search for a path linking a source node to a goal node in a huge graph is by providing the search with a sense of direction. Consider for example the problem of looking on a map for a good route between Los Angeles (LA) and New York (NY), regarding the map as a graph whose nodes are cities, and whose edges connect pairs of cities through directed routes at a cost proportional to the route length. Dijkstra’s algorithm, sketched above, provides a way for finding routes in such graphs. If our map covers all of North America, however, the algorithm would appear to behave in a strange manner, finding first shortest paths to all cities that are closer to LA than NY (like Mexico City) whether they are on the way to NY or not. The search in such a case is said to be blind, meaning that information about the goal is not used in the exploration of the map. This is definitely not the way that people search for routes in a map. If they want to go to a city that is Northeast, they will look for routes headed in that direction. The search is then focused and the total number of cities in the map is less of a problem. This basic intuition about goal-directed search was incorporated into AI search algorithms in the 1960s and 1970s in the form of *heuristic functions*: functions  $h(s)$  that provide a quick-and-dirty estimate of the distance or cost separating a state  $s$  from a goal state. In the route finding problem, for example, a useful heuristic  $h(s)$  is the *Euclidean distance* separating a city  $s$  from the target city. This is a heuristic that is very easy to compute and which provides a good approximation of the optimal cost  $h^*(s)$  of reaching the goal from  $s$ , whose exact computation would be much harder (as hard indeed as solving the original problem). Interestingly, if the evaluation function  $g(s)$  used to select the next node to label in Dijkstra’s algorithm is replaced by the more informed function  $f(s) = g(s) + h(s)$  that incorporates both the cost from  $s_0$  to  $s$  and an estimate of the cost to go from  $s$  to the goal, an *heuristic search algorithm* known as  $A^*$  is obtained that can look for paths to the goal much more effectively.<sup>23,d</sup> Heuristic search algorithms express a form of goal-driven search, where the goal is no longer passive in the search process, but actively biases the search through the heuristic term  $h(s)$ . In the extreme case in which the heuristic  $h(s)$  is perfect, that is, the estimates  $h(s)$  and the true costs  $h^*(s)$  coincide for all

states,  $A^*$  goes straight to the goal with no search at all. At the same time, if the heuristic is completely uninformative, as the null heuristic  $h(s) = 0$  for all states,  $A^*$  reduces to Dijkstra's algorithm. In the middle, when  $0 \leq h(s) \leq h^*(s)$  for all  $s$ , the heuristic  $h$  is said to be *admissible* and  $A^*$  preserves the optimality properties of Dijkstra but does less work. Indeed, the closer to  $h^*$ , the more informed the heuristic, and the less number of states that are visited in the search. Heuristic search methods have been used for finding optimal solutions to problems like Rubik's Cube that involve more than  $10^{19}$  states. The key is the use of suitable admissible heuristic functions  $h(s)$  that allow for optimal solutions while considering a tiny fraction of the problem states.<sup>25</sup>

## DOMAIN-INDEPENDENT GENERATION OF HEURISTICS

Heuristic search algorithms express a form of goal-directed search where heuristic functions are used to guide the search toward the goal. A key question is how such heuristics can be obtained for a given problem. A useful heuristic is one that provides good estimates of the cost to the goal and can be computed reasonably fast. Heuristics are traditionally devised according to the problem at hand: the Euclidean distance is a good heuristic for route finding, the sum of the Manhattan distances of each tile to its destination is a good heuristic for the sliding puzzles, and so on.<sup>e</sup> The general idea that emerges from the various heuristics developed for different problems is that heuristics  $h(s)$  can be seen as encoding the cost of reaching the goal from the state  $s$  in a problem that is simpler than the original one.<sup>24,26,27</sup> For example, the sum of the Manhattan distances in the sliding puzzles (Figure 2), corresponds to the optimal cost of a simplification of the puzzle where tiles can be moved to adjacent positions, whether such positions are empty or not. Similarly, the Euclidean heuristic for route finding is the cost of a simplification where straight routes are added between any pair of cities. The simplified problems are normally referred to as *relaxations* of the original problem.

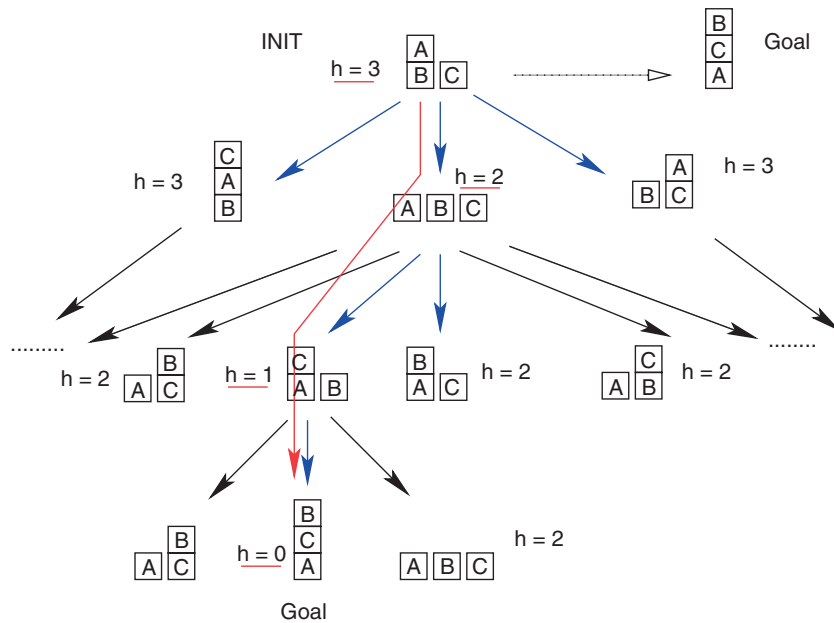
The key development in modern planning research in AI was the realization that useful heuristics could be derived automatically from the representation of the problem in a domain-independent language such as Strips.<sup>28–30</sup> It does not matter what the problem is about, a domain-independent relaxation is obtained directly and effectively from the problem representation from which heuristics that are adapted to the problem can be computed.



**FIGURE 2** | The sliding 15-puzzle where the goal is to get to a configuration where the tiles are ordered by number with the empty square last. The actions allowed are those that slide a tile into the empty square. While the problem is not simple, the heuristic that sums the horizontal and vertical distances of each tile to its target position is simple to compute and provides an informative estimate of the number of steps to the goal. In planning, heuristic functions are obtained automatically from the problem representation.

The most common and useful domain-independent relaxation in planning is the *delete-relaxation* that maps a problem  $P = \langle F, I, O, G \rangle$  in Strips into a problem  $P^+ = \langle F, I, O^+, G \rangle$  that is exactly like  $P$  but with the actions in  $O^+$  set to the actions in  $O$  with empty delete lists. That is, the delete-relaxation is a domain-independent relaxation that takes a planning problem  $P$  and produces another problem  $P^+$  where atoms are added exactly like in  $P$  but where atoms are never deleted. The relaxation implies, for example, that objects and agents can be in 'multiple places' at the same time, as when an object or an agent moves into a new place, the atom encoding the old location is not deleted. Relaxations, however, are not aimed at providing accurate models of the world; quite the opposite, simplified and even meaningless models of the world that while not accurate yield useful heuristics.

The domain-independent delete-relaxation *heuristic* is obtained as an approximation of the optimal cost of the delete-relaxation  $P^+$ , obtained from the cost of a plan that solves  $P^+$  not necessarily optimally.<sup>f</sup> The reason that an approximation is needed is because finding an optimal plan for a delete-free planning problem like  $P^+$  is still a computationally intractable task (also NP-hard). On the other hand, finding just one plan for the relaxation whether optimal or not, can be done quickly and efficiently. The property that allows for this is *decomposability*: a problem without deletes is decomposable in the sense that a plan  $\pi$  for a joint goal  $G_1$  and  $G_2$  can always be obtained from a plan  $\pi_1$  for the goal  $G_1$  and a plan  $\pi_2$  for the goal  $G_2$ . Indeed, it can be shown that the concatenation of the two plans  $\pi = \pi_1, \pi_2$  in either order, is one such plan. This property allows for a simple method for



**FIGURE 3** | A fragment of the graph corresponding to the blocks problem with the automatically derived heuristic values next to some of the nodes. The heuristic values are computed in low polynomial time and provide the search with a sense of direction. The instance can actually be solved without any search by just performing in each state the action that leads to the node with a lower heuristic value (closer to the goal). The resulting plan is shown in red; helpful actions are shown in blue.

computing plans for the relaxation from which the heuristics are derived.

The main idea behind the procedure for computing the heuristic  $h(s)$  for an arbitrary planning problem  $P$  can be explained in a few lines. For this, let  $P(s)$  refer to the problem that is like  $P$  but with the initial situation set to the state  $s$ , and let  $P^+(s)$  stand for the delete-relaxation of  $P(s)$ , that is, the problem that is like  $P(s)$  but where the delete-lists are empty. The heuristic  $h(s)$  is computed from a plan for the relaxation  $P^+(s)$  that is obtained using the decomposition property and a simple iteration. Basically, the plans for achieving the atoms  $p$  that are already true in the state  $s$ , that is,  $p \in s$ , are the empty plans, and if  $\pi_1, \pi_2, \dots, \pi_m$  are the plans for achieving each of the preconditions  $p_1, p_2, \dots, p_m$  of an action  $a$  that has the atom  $q$  in the add list,  $\pi = \pi_1, \pi_2, \dots, \pi_m$  followed by the action  $a$  is a plan for achieving  $q$ . It can be shown that this iteration yields a plan in the relaxation  $P^+(s)$  for each atom  $p$  that has a plan in the original problem  $P(s)$ , in a number of steps that is bounded by the number of variables in the problem. A plan for the actual goal  $G$  of  $P(s)$  in the relaxation  $P^+(s)$  can be obtained in a similar manner by just concatenating the plans for each of the atoms  $q$  in  $G$  in any order. Such a plan for the relaxation  $P^+(s)$ , denoted as  $\pi^+(s)$ , is usually called the *relaxed plan*. The heuristic  $h(s)$  is set to the cost of such a plan. A better estimate can be obtained if duplicate actions are removed every time plans are concatenated.<sup>31</sup> Other heuristics based also on the delete-relaxation such as the additive heuristic<sup>30</sup> and the FF heuristic,<sup>32</sup> are as informative as this heuristic

but can be computed faster. Notice that a plan  $\pi^+(s)$  for the relaxation yields the heuristic value  $h(s)$  that estimates the cost from  $s$  to the goal for the state  $s$  only. This computation thus must be repeated for every state whose heuristic value is needed.

Figure 3 displays a fragment of the directed graph corresponding to the Blocks World problem shown in Figure 1, with the automatically derived heuristic values next to some of the nodes. The heuristic values shown are computed very fast, in low polynomial time, using an algorithm similar to the one described above, with  $h(s)$  representing an approximation of the number of actions needed (cost) to solve the relaxed problem  $P^+(s)$ . Actually, the instance shown can be solved with *no search at all* by just selecting in each node, starting from the source node, the action that leads to the node with a lower heuristic value (closer to the goal). The resulting plan is shown as a red path in the figure.

In order to get a more vivid idea of where the heuristic values shown in the figure come from, consider the heuristic  $h(s)$  for the initial state where block A is on B, and both B and C are on the table. In order to get the goal ‘B on C’ in the relaxation from the state  $s$ , two actions are needed: one to get A out of the way to achieve the preconditions for moving B, the second to move B on top of C. On the other hand, in order to achieve the second goal ‘C on A’ in the relaxation from  $s$ , just the action of moving C to A is needed. The result is a heuristic value  $h(s) = 3$  as shown, which actually coincides in this case with the cost of the best plan to achieve the joint goal from  $s$  in the nonrelaxed problem  $P(s)$ . Nonetheless, this is just a coincidence, and

indeed, the best plans in the relaxation  $P^+(s)$  can be quite different than the best plans in the original problem  $P(s)$ . The best plan for  $P(s)$  is unique, moving A to the table, then C on A, and finally B on C. On the other hand, a possible optimal plan in the relaxation  $P^+(s)$  is to move first C on A, then A on the table, and finally B on C. Of course, this plan does not make any sense in the real problem where A cannot be moved when covered by C, yet the relaxation is not aimed at capturing the real problem or the real physics, it is aimed at producing informative but quick estimates of the cost to the goal. The reader can verify that for the leftmost child  $s'$  of the initial state  $s$ , the costs of the problem  $P(s')$  and the relaxation  $P^+(s')$  no longer coincide. The former is 4, while the latter is 3, the difference arising from the goal 'C on A' that in the original problem must be undone and then redone. On the other hand, in the relaxation, this is never needed as no atom is ever deleted.

One last point about the domain-independent relaxation and the resulting heuristics: the automatic computation of the heuristic value  $h(s)$  from a plan  $\pi^+(s)$  for the relaxation  $P^+(s)$  yields also relevant *structural* information that is not captured in the estimates themselves. In particular, among all the actions that are applicable in the state  $s$ , the 'relaxed plan' suggests the ones that appear as most relevant. These are the actions that are applicable in the state  $s$  and form part of the relaxed plan. In planning, such actions are said to be 'helpful',<sup>32</sup> and this structural information is used too, in one way or another in every state-of-the-art planner. Figure 3 shows the helpful actions in each of the states on the way to the goal by means of arrows depicted in blue. As it can be seen, two of the three applicable actions in the root state are helpful, just two of the six applicable actions are helpful in the second state, and just one action is helpful in the last state before the goal. It is not always the case that the best action in a state is among the ones found to be helpful in that state, or among the ones leading to the children with lowest heuristic values, yet this is often the case and this suffices for making the heuristic and the relaxation so useful. It is also interesting that the notion of helpful actions that emerges from the relaxation suggests a model for identifying the actions that appear as most promising in a state which does not require evaluating them first. This information can be used indeed for deciding which actions to consider in a state, possibly ignoring the other actions. The resulting search is not necessarily complete but it can be quite effective.<sup>32</sup> We will come back to this point below.

## HEURISTIC SEARCH

Once a heuristic  $h(s)$  is available for estimating the cost to the goal from  $s$ , the question is how to use it for finding a path to the goal. The simplest method for using the heuristic is by selecting the action  $a$  in  $s$  that produces the state  $s'$  that appears to be closest to the goal; that is, the one with minimum  $h(s')$  value. This algorithm is known as *hill-climbing*, from its use in maximization problems, as it moves the search along the slope of the heuristic function  $h$  where a value  $h(s) = 0$  denotes a goal. The problem with this and other local search algorithms is that they can get stuck in states where none of the children improves on (decreases) the heuristic value of the parent.

There are many heuristic search algorithms that avoid the problems of hill-climbing, and they can be classified into two groups: *offline* and *online* algorithms. If planning is thought of as 'thinking before acting', offline search algorithms can be understood as thinking all the way to the solution before acting, while online algorithms interleave thinking and acting. The algorithm  $A^*$  described above is an offline search algorithm with several guarantees: it is complete; that is, it will find a solution if there is one, and it is optimal if the heuristic  $h(s)$  is admissible (does not overestimate true costs). On the other hand, an algorithm like hill-climbing can be used as an online algorithm that iteratively applies the action that leads to the best child according to the heuristic. In the online setting, the algorithm is known as the *greedy algorithm*; the name greedy implying that actions are selected fast after a shallow and quick lookahead. Of course, more informed versions of this greedy online algorithm can be obtained by performing a deeper lookahead, with a lookahead of two levels, resulting in the action that leads to the best grandchild, a lookahead of three levels, resulting in the action that leads to the best great grandchild, and so on. A similar idea is used in chess playing programs that involve a slightly different model with adversarial agents.<sup>7,33</sup> Two problems with greedy algorithms enhanced with a fixed depth lookahead are that they can be trapped into a loop, returning to states that have been visited earlier in the execution, and that the computational effort grows exponentially with the lookahead depth. The loop problem in the greedy algorithm has a very elegant and powerful solution: the algorithm *Learning Real Time  $A^*$*  or *LRTA\** behaves exactly as the greedy algorithm but once it selects the action  $a$  in the state  $s$ , and before moving to the resulting state  $s'$ , it changes the heuristic value  $h(s)$  to  $c(a, s) + h(s')$ , where  $c(a, s)$  is the cost of the action  $a$  in  $s$ , and  $h(s')$  is the heuristic value of  $s'$ .<sup>34</sup> If the problem has no states from which the goal



cannot be reached (dead-ends), this simple type of learning guarantees that the goal will eventually be reached, and moreover, that if the process is restarted many times while preserving the heuristic values that have been learned, the goal will eventually be reached optimally, if the initial heuristic is admissible. The algorithm LRTA\* is also interesting because it can be easily generalized to other models where actions have probabilistic effects and states are fully or partially observable. The generalized algorithm is known as Real-Time Dynamic Programming or RTDP.<sup>35–37</sup> One last online algorithm for probabilistic models that has become popular in recent years due to its breakthrough performance in the game of Go is UCT.<sup>38,39</sup> Closely related to UCT are online variants of AO\*,<sup>40</sup> an extension of the A\* algorithm for models where actions have nondeterministic effects.<sup>24,41</sup>

## HEURISTICS, VALUES, AND APPRAISALS

Heuristic evaluation functions are used in other settings like chess playing programs and reinforcement learning. In chess, the evaluation functions are programmed by hand,<sup>24,33</sup> while in reinforcement learning, they are learned by trial-and-error.<sup>5</sup> Reinforcement learning is a family of *model-free* methods that have been shown to be effective in low-level tasks, and to provide an accurate account of learning in the brain.<sup>42</sup> Heuristic evaluation functions in domain-independent planning are computed instead using *model-based methods* where suitable relaxations are solved from scratch. These methods have been shown to work over large problems involving hundred of actions and fluents, and interesting lessons can be drawn from these methods as well. Indeed, while feelings and emotions are currently thought of as providing the appraisals that are necessary for navigating in a complex world,<sup>43</sup> there are actually very few accounts of how such appraisals may be computed. In planning, the appraisals that manage to integrate information about the current situation, the goal, and the actions, for directing the agent toward the goal, are the heuristics computed from relaxations of the problem. The computational model of appraisals that is based on the solution of relaxations that can be solved in low-polynomial time, suggests potential explanations for a number of observations. For example, the heuristics that result are ‘fast and frugal’ but unlike the heuristics considered by Gigerenzer and others,<sup>44,45</sup> they are also general: they apply to all the problems that fit the classical planning model or that can be cast in that form. The use of relaxations can also potentially

explain why appraisals may be opaque from a cognitive point of view, and thus not be conscious.<sup>46–48</sup> This is because the appraisals are obtained from a simplified model where, for example, objects can be in different places at the same time, and hence where the meaning of the symbols is different than the meaning of the symbols in the ‘true’ model. Finally, heuristic appraisals provide the agent with a sense of direction or ‘gut feeling’ that guide the action selection in the presence of many alternatives, while avoiding an infinite regress in the decision process. Indeed, the computational role of these appraisals is to avoid or at least to reduce the need to search. Explicit evaluation of all possible courses of actions is not feasible computationally, and the heuristics provide the necessary focus. Actually, as discussed above, relaxation-based heuristics in planning do not only provide an account of the value of the different options, but also of the actions that are worth evaluating; the so-called helpful actions.<sup>32</sup> This second aspect, although not necessarily in this form, may potentially help to explain a key difference between programs and humans in games such as chess, for example, where it is well known that humans consider much fewer moves than programs.

## ALTERNATIVE PLANNING METHODS

While the heuristic search approach to planning has come to dominate classical planning, many other methods have been proposed, and some of them are widely used and scale up well too. GPS, the first AI planner and one of the first AI programs, was introduced by Newell and Simon in the 50’s.<sup>1,2</sup> It introduced a technique called means-ends analysis where differences between the current state and the goal are identified and mapped into operators that can decrease those differences. Since then, the idea of means-ends analysis has been refined and extended in many ways, in the formulation of planning algorithms that are *sound* (only produce plans), *complete* (produce a plan if one exists), and *effective* (scale up to large problems). By the early 90’s, the state-of-the-art planner was UCPOP,<sup>49</sup> an implementation of an elegant planning method known as partial-order planning, where plans are not searched either forward from the starting state or backward from the goal, but are constructed from a decomposition scheme in which joint goals are decomposed into subgoals, which create as further subgoals the preconditions of the actions that can establish them.<sup>50–52</sup> The actions that are incorporated into the plan are partially ordered as needed in order to resolve possible conflicts among them. Partial-order planning algorithms are sound and complete, but do

not scale up well, as there are too many choices to make and too little guidance on how to make them.

The situation in planning changed drastically in the middle 90's with the introduction of Graphplan,<sup>53</sup> an algorithm that appeared to have little in common with previous approaches but scaled up much better. Graphplan builds a *plan graph* in polynomial time reasoning forward from the initial state, which it then searches backward from the goal to find a plan. It was shown later that the reason Graphplan scaled up well was due to a powerful admissible heuristic implicit the plan graph.<sup>54</sup> The success of Graphplan prompted other approaches. In the SAT approach,<sup>55</sup> the planning problem for a fixed planning horizon is converted into a general *satisfiability* problem expressed as a set of clauses (a formula in conjunctive normal form or CNF) that is fed into state-of-the-art SAT solvers that currently manage to solve huge SAT instances even if the SAT problem is NP-complete.<sup>56</sup> A clause is a disjunction of literals, that is, propositional symbols or their negations as in  $x \vee \neg y \vee z$ , and a set of clauses is satisfiable if there is a truth assignment to the symbols such that each clause has at least one true literal. In the CNF encoding of the planning problem, atoms and actions are indexed with time indices that range from zero until a planning horizon that is increased one by one until the resulting clauses are satisfiable and a plan can be read from the satisfying assignment. Due to the excellent performance of current SAT solvers, SAT planners manage to solve large problems very fast, making them practically competitive with heuristic search planners.<sup>57</sup> State-of-the-art heuristic search planners use heuristic values derived from the delete-relaxation,<sup>28,29</sup> information about the action that are most helpful,<sup>32</sup> and implicit subgoals of the problem, called landmarks, that are also extracted automatically from the problem with methods similar to those used for deriving heuristics.<sup>58,59</sup> Recently, the success of classical planners has also been explained in term of a structural width parameter that appears to be bounded and small in many domains when goals are restricted to single atoms. Such problems can be solved in time that is exponential in their width, while problems with joint goals can often be decomposed easily into a sequence of low-width problems with single goals.<sup>60</sup>

## RICHER PLANNING MODELS

Classical planning is planning with deterministic actions from an initial state that is fully known. Many planning problems, however, involve features that are not part of this basic model such as uncertainty,

incomplete information, and soft goals. Two types of methods have been pursued for dealing with such features: a *top-down* approach, where *native solvers* have been developed for more expressive planning models, and a *bottom-up* approach, where the power of classical planners is exploited by means of *translations*.

## MDP and POMDP Planning

MDP and POMDP planners are examples of native solvers for more expressive planning models that accommodate stochastic actions, and either full or partial state observability. A Markov Decision Process (MDP) is a state model where the state transition function  $f(a, s)$  is replaced by state transition probabilities  $P_a(s'|s)$ , and the next state  $s'$ , while no longer predictable with certainty, is assumed to be *fully observable*.<sup>61-63</sup> A solution to an MDP is a function  $\pi$  from states into actions, called a *policy*, that drives the system to a goal state with certainty. A policy induces a probability distribution over the possible state trajectories, and since every state trajectory has a cost, an optimal policy is a policy that drives the system to the goal at a minimum expected cost.

Partially Observable MDPs (POMDPs) extend MDPs by relaxing the assumptions that states are fully observable.<sup>61,64,65</sup> In a POMDP, a set of observation tokens  $o \in O$  is assumed along with a *sensor model*  $P_a(o|s)$  that relates the true but hidden state  $s$  of the system with the observable token  $o$ . In POMDPs, the initial state of the system is not known but is characterized by a probability distribution, and the task is to drive the system to a final, fully observable target state. Solutions to POMDPs are closed-loop controllers that map *belief states* into actions, with optimal solutions reaching the target state at a minimum expected cost. The belief states are probability distributions over the states.

A simple example of a POMDP is a robot that has to reach a certain location by means of actions whose effects can only be predicted probabilistically. The state of the problem, which is the location of the robot, is not fully observable but rather sonars or other feedback mechanisms are used to provide the robot with partial information about the state. A map showing the locations that are blocked by obstacles may be available, but the map cannot be used in a direct way, as the robot does not know with certainty its location in the map. On the other hand, if the robot can have perfect access to its location, the problem is not a POMDP but an MDP.

While MDPs and POMDPs are commonly described using positive or negative *rewards* rather

than positive *costs*, and using *discount factors* rather than *goal states*, simple transformations are known for translating discounted reward MDPs and POMDPs into equivalent *goal MDPs* and *goal POMDPs* as above that are strictly more expressive.<sup>37,61</sup> From a computational point of view, traditional dynamic programming algorithms have been used to solve MDPs and POMDPs offline,<sup>61,66,67</sup> yet the methods that scale up best to larger problems are online methods, closely related to the online heuristic search algorithms considered above for deterministic problems. These include Real-time Dynamic Programming,<sup>35–37,68</sup> and UCT.<sup>38,40,69,70</sup> These algorithms are also related to reinforcement learning methods, which can be characterized as methods for solving an MDP by trial-and-error without knowing the value of the cost and probability parameters.<sup>5</sup> Normally, reinforcement learning algorithms learn a representation of the optimal MDP policy without learning the value of these parameters. These are called *model-free* MDP methods.<sup>71,72</sup> On the other hand, there are approaches that incrementally learn the MDP model and derive the policy from the model. These are called *model-based* reinforcement learning methods, and the most effective of them reduce the learning problem to a planning problem over an ‘optimistic’ model that is refined incrementally.<sup>73–75</sup> A last class of POMDP methods map probabilistic planning problems into probabilistic reasoning problems over Bayesian Networks,<sup>76–78</sup> in analogy to the SAT approach to deterministic planning, where atoms and actions are indexed in time up to a given planning horizon. While the reduction of planning to inference, deductive or probabilistic, is appealing, the scalability-quality tradeoff in the probabilistic case, unlike the SAT approach in the deductive case, is not yet clear in comparison with state-of-the-art methods.

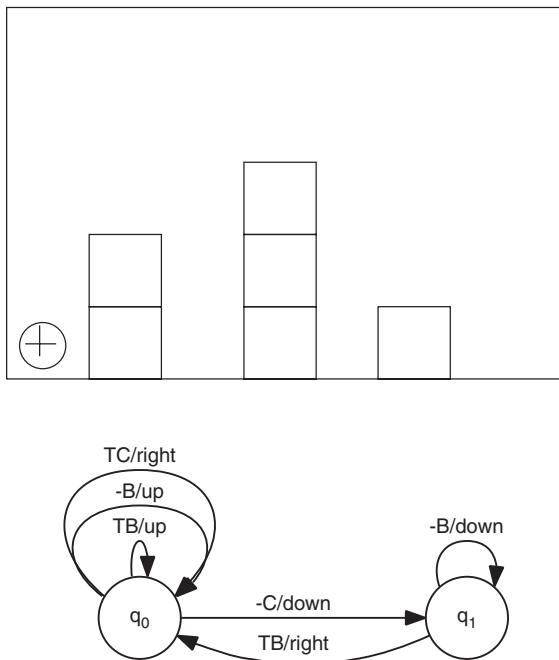
## Translations into Classical Planning

Translation-based approaches handle features that are absent from the classical model such as soft-goals, plan-constraints, extended temporal goals, uncertainty, and partial feedback, by compiling them away. We focus only on translations that preserve the semantics of the original problem, yet approximate translations have been found useful as well. For example, FF-Replan is an online planner for MDPs whose value has been proved in some of the MDP planning competitions held so far.<sup>79</sup> FF-Replan’s approach to MDP planning is simple but effective. It selects the action to do next using a simple relaxation of the MDP: first the probabilities associated with the nondeterministic effects of a probabilistic action

are dropped, then the nondeterministic effects of the same action are mapped into different deterministic actions. The relaxation thus maps an MDP where the uncertain effects are controlled by nature into a deterministic problem where the uncertain effects are controlled by the planning agent. Once the plan is obtained for the relaxation, the plan is executed until the agent finds itself in a state that is different from the one predicted by the relaxation. The process resumes from this state and the execution terminates when the goal is reached. The resulting online planner is not optimal and can get trapped into dead-ends, yet in most domains, the simplification works rather well. A similar idea, where closed-loop feedback controls for stochastic systems are designed using a simplified model, is common in control engineering.<sup>80</sup>

Features that are absent from the classical planning model include rewards, and in particular *soft goals*. Soft goals refer to formulas that if achieved along with the goal entail a positive utility. In the presence of soft goals, the task is to compute action sequences that achieve the goal and maximize overall utility, rather than sequences that achieve the goal and minimize cost. The utility is defined as the sum of the utilities of the soft goals achieved minus the cost of the action sequence. Soft goals express soft preferences as opposed to hard goals that express hard preferences. For example, getting a new T-shirt can be a hard goal for John, while getting a second T-shirt can be regarded as a soft goal. All plans must thus deliver John a new T-shirt, but whether the best plans will deliver him a second T-shirt or not will depend on the extra costs and utilities involved. Interestingly, it turns out that soft-goals can be compiled away easily and efficiently resulting into standard classical planning problems with hard goals only.<sup>81</sup> For soft-goal atoms  $A$  one just needs to create new atoms  $A'$  and make them into hard goals that can be achieved in one of two ways: by means of the new action  $\text{collect}(A)$  with precondition  $A$  and cost 0, or by means of the new action  $\text{forego}(A)$  with no precondition and cost equal to the utility of  $A$ . The best plans for the soft goals become the best plans for the resulting classical problem.

*Uncertain information* can also be compiled away in conformant problems with deterministic actions.<sup>82</sup> These are problems where the initial state is partially known and a plan to the goal is sought that would work for any possible initial state. The translation maps a conformant problem  $P$  into a classical problem, whose solutions, computable with off-the-shelf classical planners, encode the solutions to  $P$ . The complexity of the translation is exponential in a width parameter that is often bounded and small. The



**FIGURE 4** | Top: Problem where visual-marker (circle on the lower left) must be placed on top of a green block by just observing what's on the marked cell. Down: Finite-state controller obtained with a classical planner from suitable translation. The controller has two states, the initial state  $q_0$  and the state  $q_1$ . An edge  $q \rightarrow q'$  with label  $o/a$  means to do the action  $a$  when observing  $o$  in the state  $q$ , and to move then to the state  $q'$ . The controller shown solves the problem, and any variation of it resulting from changes in the number or configuration of blocks.

ideas behind this translation have been used since to define effective action selection mechanisms for online planning in the presence of *partial observability*<sup>83–85</sup> and for computing solutions to planning problems with uncertainty and partial feedback in the form of *finite-state controllers*. One such finite-state controller is shown in Figure 4 for a problem inspired by the use of deictic representations where objects are not assumed to have unique names but can be referred to by means of suitable indexical expressions and markers.<sup>86,87</sup> In the figure, a visual marker, that is the circle on the lower left, must be placed on top of a green block by moving it one cell at a time. The location of the green block is not known (it can be anywhere), and the observations are whether the cell currently marked contains a green block (G), a nongreen block (B), or neither (C), and also whether the marked cell is at the level of the table (T) or not (–). The finite-state controller shown below has been obtained by running a *classical planner* on a suitable translation of the problem.<sup>88</sup> The controller has two states, the initial state  $q_0$  and the state  $q_1$ . An edge  $q \rightarrow q'$  with label  $o/a$  means to do the action  $a$  when observing  $o$  in the state  $q$ , and to move then to the state

$q'$ . For example, when the observation TC is received in  $q_0$ , the action of moving the marker to the Right is done, and the controller remains in the same state. The reader can check that the controller obtained with the classical planner not only solves the original problem on the left, but also any modification of it resulting from changes in the number or configuration of blocks. In other words, due to its representation, the solution to the original problem generalizes to a much larger class of problems. The problem of devising controllers or strategies that work for many or even all domain instances is usually referred to as *generalized planning*.<sup>89,90</sup>

## CHALLENGES

There are several open problems in planning research. One of them is planning in the presence of other agents that plan, often called *multiagent planning*. People do this naturally all the time: walking on the street, driving, etc. The first question is how plans should be defined. This is a subtle problem and many proposals have been put forward, often building on equilibria notions from game theory.<sup>91–93</sup> Yet, there are currently no models, algorithms, or implementations of domain-independent planners able to plan meaningfully and efficiently in such settings. This is not entirely surprising, however, given the known limitations of game theory as a descriptive theory of human behavior.<sup>94</sup> Eventually, a working theory of multiagent planning could shed light on the computation, nature, and role of social emotions in multiagent settings, very much as single-agent planning may shed light on the computation, nature, and role of goal appraisals and heuristics in the single-agent setting. A second open problem is the automatic construction and use of hierarchies. Hierarchies form a basic component of Hierarchical Task Networks or HTNs, an alternative model for planning that is concerned with the encoding of *strategies for solving problems*,<sup>8,51,95</sup> yet hierarchies play no role in state-of-the-art domain independent planners which are completely flat. There is a large body of work on abstract problem solving that is relevant to this question, both old<sup>96–99</sup> and new,<sup>100–102</sup> but no robust answer yet. A third open problem is learning the planning models by interacting with the environment. We have discussed model-based reinforcement learning algorithms that actively learn model parameters such as probabilities and rewards,<sup>73–75</sup> yet a harder problem is learning the states themselves from partial observations. Several attempts to generalize reinforcement learning algorithms to partially observable settings have been made,

some of which learn to identify useful features and feature histories,<sup>103–105</sup> but none so far that can come up with the states and models themselves in a robust and scalable manner.

## DISCUSSION

The relevance of the early work in AI to Cognitive Science was based on *intuition*: programs provided a way for specifying intuitions precisely and for trying them out. The more recent work on *domain-independent solvers* in AI is more technical and experimental, and is focused not on reproducing intuitions but on *scalability*. This may give the impression that recent work in AI is less relevant to Cognitive Science than work in the past. This impression, however, may prove wrong on at least two grounds. First, intuition is not what it used to be, and it is now regarded as the tip of an iceberg whose bulk is made of massive amounts of shallow, fast, but unconscious inference mechanisms that cannot be rendered explicit.<sup>45–47</sup> Second, whatever these mechanisms are, they appear to work pretty well and to scale up. This is no small feat, given that most methods, whether intuitive or not, do not. By focusing on the study of meaningful models and the computational methods for dealing with them *effectively*, AI may prove its relevance to the understanding of human cognition in ways that may go well beyond the rules, cognitive architectures, and knowledge structures of the 80's. Human cognition, indeed, still provides the inspiration and motivation for a lot of research in AI. The use of Bayesian networks in developmental psychology for understanding how children acquire and use causal relations,<sup>106</sup> and the use of reinforcement learning algorithms in neuroscience for interpreting the activity of dopamine cells in the brain,<sup>42</sup> are two examples of general AI techniques that have made it recently into cognitive science. As AI focuses on models and solvers able to scale up, more techniques are likely to follow. In this paper, we have reviewed work in computational

models of planning in AI, with an emphasis on deterministic planning models where automatically derived relaxations and heuristics manage to integrate information about the current situation, the goal, and the actions for directing the agent effectively toward the goal. The computational model of goal appraisals that is based on the solution of low-polynomial relaxations may shed light on the computation, nature, and role of other types of appraisals, and on why appraisals are opaque to cognition and cannot be rendered conscious or articulated in words.

## NOTES

<sup>a</sup>More precisely, in the language of computer science, planning is said to be NP-hard, meaning that there cannot be a general sound and complete planning algorithm that runs in polynomial time, unless the fundamental conjecture in computer science that NP-hard problems do not admit polynomial solutions, widely believed to be true, turns out to be false.<sup>11</sup>

<sup>b</sup>The relationship between planning and emotion has been considered by other authors that aim to import elements from the theory of emotions into AI planners.<sup>15,16</sup> Our focus here goes in the opposite direction, on what modern computational models of planning can tell us about the generation, nature, and role of emotional appraisals.

<sup>c</sup>The age of the universe is estimated at  $13.7 \times 10^9$  years approximately. Generating  $2^{100}$  nodes at  $10^7$  nodes a second would take in the order of  $10^{15}$  years, as  $2^{100}/(10^7 \times 60 \times 60 \times 24 \times 365) \approx 4 \times 10^{15}$ .

<sup>d</sup>This description of A\* assumes that the heuristic is monotonic. A more complete description of A\* and variations can be found in the standard textbooks.<sup>7,24</sup>

<sup>e</sup>The Manhattan distance of a tile is the sum of the horizontal and vertical distances that separate its current position from its destination.

<sup>f</sup>Heuristics derived in this way are not admissible; that is, they may overestimate the true costs and hence are not suitable for computing optimal plans.

## ACKNOWLEDGMENTS

The author thanks the reviewers for useful comments. H. Geffner is partially supported by grants TIN2009-10232 and CSD2010-00034 (SimulPast), MICINN, Spain, and EC-7PM-SpaceBook.

## REFERENCES

1. Newell A, Simon H. Elements of a theory of human problem solving. *Psychol Rev* 1958, 65:151–166.
2. Newell A, Simon H. GPS: a program that simulates human thought. In: Feigenbaum E, Feldman J, eds. *Computers and Thought*. McGraw Hill; 1963, 279–293.
3. Brooks R. A robust layered control system for a mobile robot. *IEEE J Robot Autom* 1987, 2:14–27.

4. Mataric MJ. *The Robotics Primer*. Cambridge, MA: MIT Press; 2007.
5. Sutton R, Barto A. *Introduction to Reinforcement Learning*. Cambridge, MA: MIT Press; 1998.
6. Mitchell T. *Machine Learning*. Boston, MA: McGraw-Hill; 1997.
7. Russell S, Norvig P. *Artificial Intelligence: A Modern Approach*. 3rd ed. Upper Saddle River, NJ: Prentice Hall; 2009.
8. Ghallab M, Nau D, Traverso P. *Automated Planning: Theory and Practice*. San Francisco, CA: Morgan Kaufmann; 2004.
9. Geffner H, Bonet B. *Advanced Introduction to Planning: Models and Methods*. San Rafael, CA: Morgan & Claypool; 2013.
10. Dennett D. *Kinds of Minds*. New York: Basic Books; 1996.
11. Bylander T. The computational complexity of STRIPS planning. *Artif Intell* 1994, 69:165–204.
12. Tooby J, Cosmides L. The psychological foundations of culture. In: Barkow J, Cosmides L, Tooby J, eds. *The Adapted Mind*. Oxford, NY: Oxford University Press; 1992.
13. Coles AJ, Coles A, Olaya AG, Jiménez S, López CL, Sanner S, Yoon S. A survey of the seventh international planning competition. *AI Mag* 2012, 33:83–88.
14. Geffner H. Heuristics, planning, cognition. In: Dechter R, Geffner H, Halpern J, eds. *Heuristics, Probability and Causality. A Tribute to Judea Pearl*. London: College Publications; 2010.
15. Gratch J. Why you should buy an emotional planner. *Proceedings of Agents'99 Workshop on Emotion-based Agent Architectures (EBAA'99)*, 1999. Available at: <http://emotions.usc.edu/~gratch/gratch-ebaa99.pdf>.
16. Gratch J, Marsella S. A domain-independent framework for modeling emotion. *Cogn Syst Res* 2004, 5:269–306.
17. Fikes R, Nilsson N. STRIPS: A new approach to the application of theorem proving to problem solving. *Artif Intell* 1971, 1:27–120.
18. McDermott D. PDDL—the planning domain definition language, 1998. Available at: <http://ftp.cs.yale.edu/pub/mcdermott>
19. Backstrom C, Nebel B. Complexity results for SAS+ planning. *Comput Intell* 1995, 11:625–655.
20. Younes H, Littman M, Weissman D, Asmuth J. The first probabilistic track of the international planning competition. *J Artif Intell Res* 2005, 24:851–887.
21. Dijkstra E. A note on two problems in connexion with graphs. *Numer Math* 1959;1:269–271.
22. Cormen TH, Leiserson CE, Rivest RL, Stein C. *Introduction to Algorithms*. Cambridge, MA: The MIT Press; 2009.
23. Hart P, Nilsson N, Raphael B. A formal basis for the heuristic determination of minimum cost paths. *IEEE Trans Syst Sci Cybern* 1968, 4:100–107.
24. Pearl J. *Heuristics*. Reading, MA: Addison Wesley; 1983.
25. Korf R. Finding optimal solutions to Rubik's cube using pattern databases. *Proceedings of AAAI-98*, 1998, 1202–1207.
26. Simon H. A behavioral model of rational choice. *Quart J Econ* 1955, 69:99–118.
27. Minsky M. Steps toward artificial intelligence. *Proc IRE* 1961, 49:8–30.
28. McDermott D. A heuristic estimator for means-ends analysis in planning. *Proceedings of AIPS-96*, 1996, 142–149.
29. Bonet B, Loerincs G, Geffner H. A robust and fast action selection mechanism for planning. *Proceedings of AAAI-97*, 1997, 714–719.
30. Bonet B, Geffner H. Planning as heuristic search. *Artif Intell* 2001, 129:5–33.
31. Keyder E, Geffner H. Heuristics for planning with action costs revisited. *Proceedings of ECAI-08*, 2008, 588–592.
32. Hoffmann J, Nebel B. The FF planning system: fast plan generation through heuristic search. *J Artif Intell Res* 2001, 14:253–302.
33. Newell A, Shaw JC, Simon H. Chess-playing programs and the problem of complexity. In: Feigenbaum E, Feldman J, eds. *Computers and Thought*. McGraw Hill; 1963, 109–133.
34. Korf R. Real-time heuristic search. *Artif Intell* 1990, 42:189–211.
35. Barto A, Bradtke S, Singh S. Learning to act using real-time dynamic programming. *Artif Intell* 1995, 72:81–138.
36. Bonet B, Geffner H. Labeled RTDP: Improving the convergence of real-time dynamic programming. *Proceedings of 13th International Conference on Automated Planning and Scheduling (ICAPS-2003)*. Menlo Park, CA: AAAI Press; 2003, 12–31.
37. Bonet B, Geffner H. Solving POMDPs: RTDP-Bel vs. point-based algorithms. *Proceedings of IJCAI*, 2009, 1641–1646.
38. Kocsis L, Szepesvári C. Bandit based Monte-Carlo planning. *Proceedings of ECML-2006*, 2006, 282–293.
39. Gelly S, Silver D. Combining online and offline knowledge in UCT. *Proceedings of ICML*, 2007, 273–280.
40. Bonet B, Geffner H. Action selection for MDPs: Anytime AO\* vs. UCT. *Proceedings of AAAI-2012*, 2012.
41. Martelli A, Montanari U. Additive AND/OR graphs. *Proceedings of IJCAI-73*, 1973, 1–11.
42. Schultz W, Dayan P, Montague P. A neural substrate of prediction and reward. *Science* 1997, 275:1593–1599.

43. Damasio A. *Descartes' Error: Emotion, Reason, and the Human Brain*. Castleton, NY: Quill; 1995.
44. Gigerenzer G, Todd P. *Simple Heuristics that Make Us Smart*. Oxford, NY: Oxford University Press; 1999.
45. Gigerenzer G. *Gut Feelings: The Intelligence of the Unconscious*. New York: Viking Books; 2007.
46. Wilson T. *Strangers to Ourselves*. Cambridge, MA: Belknap Press; 2002.
47. Hassin R, Uleman J, Bargh J. *The New Unconscious*. Oxford, NY: Oxford University Press; 2005.
48. Kahneman D. *Thinking, fast and slow*, Farrar, Straus and Giroux, New York; 2011.
49. Penberthy J, Weld D. UCPOP: a sound, complete, partial order planner for ADL, *Proceedings of KR-92*, 1992.
50. Sacerdoti E. The nonlinear nature of plans. *Proceedings of IJCAI-75*, 1975, 206–214.
51. Tate A. Generating project networks. *Proceedings of IJCAI*, 1977, 888–893.
52. McAllester D, Rosenblitt D. Systematic nonlinear planning. *Proceedings of AAAI-91*, 1991, 634–639.
53. Blum A, Furst M. Fast planning through planning graph analysis. *Proceedings of IJCAI-95*. 1995, 1636–1642.
54. Haslum P, Geffner H. Admissible heuristics for optimal planning. *Proceedings of the Fifth International Conference on AI Planning Systems (AIPS-2000)*, 2000, 70–82.
55. Kautz H, Selman B. Pushing the envelope: planning, propositional logic, and stochastic search. *Proceedings of AAAI*, 1996, 1194–1201.
56. Biere A, Heule M, van Maaren H, Walsh T. *Handbook of Satisfiability: Volume 185 Frontiers in Artificial Intelligence and Applications*. Amsterdam: IOS Press; 2009.
57. Rintanen J. Heuristics for planning with SAT. *Proceedings of on Principles and Practice of Constraint Programming (CP 2010)*. Berlin: Springer; 2010, 414–428.
58. Hoffmann J, Porteous J, Sebastia L. Ordered landmarks in planning. *J Artif Intell Res* 2004, 22:215–278.
59. Richter S, Westphal M. The LAMA planner: guiding cost-based anytime planning with landmarks. *J Artif Intell Res* 2010, 39:127–177.
60. Lipovetzky N, Geffner H. Width and serialization of classical planning problems. *Proceedings of ECAI*. Amsterdam: IOS Press; 2012, 540–545.
61. Bertsekas D. *Dynamic Programming and Optimal Control, Vols 1 and 2*. Nashua, NH: Athena Scientific; 1995.
62. Puterman M. *Markov Decision Processes –Discrete Stochastic Dynamic Programming*. New York: John Wiley and Sons, Inc.; 1994.
63. Boutilier C, Dean T, Hanks S. Decision-theoretic planning: structural assumptions and computational leverage. *J Artif Intell Res (JAIR)* 1999, 11:1–94.
64. Astrom K. Optimal control of Markov decision processes with incomplete state estimation. *J Math Anal Appl* 1965, 10:174–205.
65. Kaelbling L, Littman M, Cassandra T. Planning and acting in partially observable stochastic domains. *Artif Intell* 1998, 101:99–134.
66. Bellman R. *Dynamic Programming*. Princeton, NJ: Princeton University Press; 1957.
67. Howard R. *Dynamic Probabilistic Systems–Volume I: Markov Models*. New York: John Wiley and Sons; 1971.
68. Kolobov A, Mausam, Weld D. LRTDP vs. UCT for online probabilistic planning. *Proceedings of AAAI*, 2012.
69. Silver D, Veness J. Monte-Carlo planning in large POMDPs. *Advances in Neural Information Processing Systems (NIPS)*, 2010. 2164–2172.
70. Keller T, Eyerich P. PROST: Probabilistic planning based on UCT. *Proceedings of ICAPS*, 2012.
71. Watkins C, Dayan P. Q-learning. *Mach Learn* 1992, 8:279–292.
72. Sutton R. Learning to predict by the methods of temporal differences. *Mach Learn* 1988, 3:9–44.
73. Kearns M, Singh S. Near-optimal reinforcement learning in polynomial time. *Mach Learn* 2002, 49:209–232.
74. Brafman R, Tenenbholz M. R-Max: a general polynomial time algorithm for near-optimal reinforcement learning. *J Mach Learn Res* 2003, 3:213–231.
75. Asmuth J, Littman M. Learning is planning: near bayes-optimal reinforcement learning via Monte-Carlo tree search. *Proceedings of UAI*, 2011, 19–26.
76. Attias H. Planning by probabilistic inference. *Proceedings of the 9th International Workshop on Artificial Intelligence and Statistics*, 2003.
77. Toussaint M, Storkey A. Probabilistic inference for solving discrete and continuous state Markov decision processes. *Proceedings of 23rd International Conference on Machine Learning*, 2006, 945–952.
78. Botvinick M, An J. Goal-directed decision making in the prefrontal cortex: a computational framework. *Adv Neural Inform Process Syst (NIPS)*, 2008. 169–176.
79. Yoon S, Fern A, Givan R. FF-replan: a baseline for probabilistic planning. *Proceedings of ICAPS-07*, 2007, 352–359.
80. Ogata K. *Modern Control Engineering*. Upper Saddle River, NJ: Prentice Hall; 2001.
81. Keyder E, Geffner H. Soft goals can be compiled away. *J Artif Intell Res* 2009, 36:547–556.

82. Palacios H, Geffner H. Compiling uncertainty away in conformant planning problems with bounded width. *J Artif Intell Res* 2009, 35:623–675.
83. Albore A, Palacios H, Geffner H. A translation-based approach to contingent planning. *Proceedings of IJCAI-09*, 2009, 1623–1628.
84. Bonet B, Geffner H. Planning under partial observability by classical replanning: theory and experiments. *Proceedings of IJCAI-11*, 2011, 1936–1941.
85. Shani G, Brafman R. Replanning in domains with partial information and sensing actions. *Proceedings of IJCAI-2011*, 2011, 2021–2026.
86. Chapman D. Penguins can make cake. *AI Mag* 1989, 10:45–50.
87. Ballard D, Hayhoe M, Pook P, Rao R. Deictic codes for the embodiment of cognition. *Behav Brain Sci* 1997, 20:723–742.
88. Bonet B, Palacios H, Geffner H. Automatic derivation of memoryless policies and finite-state controllers using classical planners. *Proceedings of ICAPS-09*, 2009, 34–41.
89. Martin M, Geffner H. Learning generalized policies from planning examples using concept languages. *Appl Intell* 2004, 20:9–19.
90. Hu Y, De Giacomo G. Generalized planning: synthesizing plans that work for multiple environments. *Proceedings of IJCAI*, 2011; 918–923.
91. Bowling M, Jensen R, Veloso M. A formalization of equilibria for multiagent planning. *Proceedings of IJCAI-03*, 2003, 1460–1462.
92. Larbi R, Konieczny S, Marquis P. Extending classical planning to the multi-agent case: a game-theoretic approach. *Proceedings of 9th European Conf. on Symbolic and Quantitative Approaches to Reasoning with Uncertainty (ECSQARU 2007)*, volume 4724 of *Lecture Notes in Computer Science*. Springer; 2007, 731–742.
93. Brafman R, Domshlak C, Engel Y, Tennenholtz M. Planning games. *Proceedings of IJCAI-09*, 2009, 73–78.
94. Camerer C. *Behavioral Game Theory: Experiments in Strategic Interaction*. Princeton, NJ: Princeton University Press; 2003.
95. Erol K, Hendler J, Nau DS. HTN planning: Complexity and expressivity. *Proceedings of AAAI-94*, 1994, 1123.
96. Sacerdoti E. Planning in a hierarchy of abstraction spaces. *Artif Intell* 1974, 5:115–135.
97. Korf R. Planning as search: a quantitative approach. *Artif Intell* 1987, 33:65–88.
98. Knoblock C. Learning abstraction hierarchies for problem solving. *Proceedings of AAAI-90*, 1990, 923–928.
99. Bacchus F, Yang Q. Downward refinement and the efficiency of hierarchical problem solving. *Artif Intell* 1994, 71:43–100.
100. McIlraith S, Fadel R. Planning with complex actions. *Proceedings of NMR-02*, 2002, 356–364.
101. Marthi B, Russell S, Wolfe J. Angelic semantics for high-level actions. *Proceedings of ICAPS-07*, 2007, 232–239.
102. Jonsson A. The role of macros in tractable planning over causal graphs. *Proceedings of IJCAI-07*, 2007, 1936–1941.
103. McCallum A. Overcoming incomplete perception with utile distinction memory. *Proceedings Tenth International Conference on Machine Learning*, 1993, 190–196.
104. Stanley K, Bryant B, Miikkulainen R. Real-time neuroevolution in the nero video game. *IEEE Trans Evolution Comput* 2005, 9:653–668.
105. Veness J, Ng K, Hutter M, Uther W, Silver D. A Monte-Carlo AIXI approximation. *J Artif Intell Res* 2011, 40:95–142.
106. Gopnik A, Glymour C, Sobel D, Schulz L, Kushnir T, Danks D. A theory of causal learning in children: causal maps and Bayes nets. *Psychol Rev* 2004, 111:3–31.