



DIGITAL ACCESS TO SCHOLARSHIP AT HARVARD

Plan Recognition in Exploratory Domains

The Harvard community has made this article openly available.
[Please share](#) how this access benefits you. Your story matters.

| | |
|--------------------------|--|
| Citation | Gal, Ya'akov, Swapna Reddy, Stuart M. Shieber, Andee Rubin, and Barbara J. Grosz. 2012. Plan recognition in exploratory domains. <i>Artificial Intelligence</i> 176(1): 2270-2290. |
| Published Version | doi:10.1016/j.artint.2011.09.002 |
| Accessed | July 25, 2017 10:47:01 PM EDT |
| Citable Link | http://nrs.harvard.edu/urn-3:HUL.InstRepos:5343166 |
| Terms of Use | This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP |

(Article begins on next page)

Plan Recognition in Exploratory Domains

Ya'akov Gal^{1,2}, Swapna Reddy², Stuart Shieber², Andee Rubin³, and Barbara Grosz²

¹Department of Information Systems Engineering, Ben-Gurion University of the Negev, Israel

²School of Engineering and Applied Sciences, Harvard University, USA

³TERC, USA

Abstract

This paper describes a challenging plan recognition problem that arises in environments in which agents engage widely in exploratory behavior, and presents new algorithms for effective plan recognition in such settings. In exploratory domains, agents' actions map onto logs of behavior that include switching between activities, extraneous actions, and mistakes. Flexible pedagogical software, such as the application considered in this paper for statistics education, is a paradigmatic example of such domains, but many other settings exhibit similar characteristics. The paper establishes the task of plan recognition in exploratory domains to be NP-hard and compares several approaches for recognizing plans in these domains, including new heuristic methods that vary the extent to which they employ backtracking, as well as a reduction to constraint-satisfaction problems. The algorithms were empirically evaluated on people's interaction with flexible, open-ended statistics education software used in schools. Data was collected from adults using the software in a lab setting as well as middle school students using the software in the classroom. The constraint satisfaction approaches were complete, but were an order of magnitude slower than the heuristic approaches. In addition, the heuristic approaches were able to perform within 4% of the constraint satisfaction

approaches on student data from the classroom, which reflects the intended user population of the software. These results demonstrate that the heuristic approaches offer a good balance between performance and computation time when recognizing people's activities in the pedagogical domain of interest.

1 Introduction

In this paper we report on the development and evaluation of algorithms for recognizing users' plans in domains in which users engage in exploratory and error-prone behaviors. The challenges presented by these domains were made evident by our work with students using open-ended computer software for learning statistics, but they arise in human-computer interaction more broadly.

Indeed, developing technology is changing rote and monolithic interaction styles between computers and their users to more flexible types of interactions that allow users to explore and interleave between different activities. Examples of these flexible systems include interactive drawing tools (Ryall et al., 1997), Integrated Development Environments (IDEs), collaborative writing assistants (Babaian et al., 2002), computer games, and educational software (Yaron et al., 2009).

To be effective partners, these systems need to recognize the activities their users are carrying out and to use that information to provide support in a way that guides users' interactions effectively. For example, an intelligent drawing tool may infer that several objects on the canvas are all representatives of the same class. When the user modifies an attribute in one of the forms, the system will identify and duplicate this change in the other objects in the class. Another benefit of recognizing users' activities in software is to provide assessments of user performance. Such capabilities in educational and pedagogical systems could increase teachers' abilities to identify those students who are having difficulty.

Classical approaches to plan recognition have assumed a goal-oriented agent whose activities are consistent with the recognizers' knowledge base and who forms a single encompassing plan. In contrast, flexible systems allow users to follow multiple plans, interleave actions from different plans, and perform redundant actions; they also tolerate user mistakes. Thus, inferring users' plans in these systems gives rise to a more complex sort of plan

recognition problem.

This paper presents several new algorithms for keyhole plan recognition in exploratory domains¹. The algorithms are post-hoc, in that they infer plans from complete interaction sequences, rather than after each observed action, as in on-line recognition (Bui, 2003). The algorithms we present vary in completeness (that is, whether plans are guaranteed to be found) and computational complexity. We investigate the trade-off between completeness and complexity empirically, by comparing the performance of different plan recognition algorithms on real-world data.

Our empirical analysis uses an educational software system for statistics education. Educational software is increasingly designed to be open-ended and flexible in order to support the types of exploratory activities that facilitate students' learning experience. This gives students the resources to explore concepts in new ways, but their interactions may be erratic or unfocused, making it challenging to recognize plans. During the chaos of a lab session, it is impossible for teachers to track each student's progress. As a result it is difficult to adapt their teaching to their students' work. Educational software thus provides an important domain for plan recognition. A well structured post-hoc representation of the plans behind students' activities would enable teachers to make better pedagogical decisions in the classroom.

The research we report used a commercial system called TinkerPlots, used world-wide to teach students in grades 4 through 8 about statistics and mathematics (Konold and Miller, 2004). Using TinkerPlots, students build stochastic models and generate pseudo-random samples to analyze the underlying probability distributions. Our study used four different problems for which students interacted with TinkerPlots to model hypothetical situations and to determine the probability of events.

Students' interactions with TinkerPlots are complex. They may pursue multiple plans and interleave actions from different plans. They may be confused about the appropriate plan to take, and they may make mistakes. These behaviors create a challenging domain for plan recognition algorithms. Any number of extraneous actions may be interleaved among those that are a part of a successful plan. In addition, actions that are crucial to successful plans may occur in almost any order.

¹We use the term "keyhole plan recognition," coined by Cohen et al. (1981), to refer to the fact that the acting agent is not signalling its plan to the observer.

All of the algorithms presented in the paper compose (possibly non-contiguous) interaction sequences from users' interactions into a series of interdependent tasks and sub-tasks. They infer students' plans by comparing their interaction sequence to ideal solutions, or recipes, that were specified by domain experts. At the end of this process, the algorithms output a hierarchical plan that explains the student's strategy during the session. The algorithms separate those actions that contribute to solving the problem from extraneous actions and mistakes.

This paper integrates and extends initial reports of past studies (Gal et al., 2008; Reddy et al., 2009) and makes several contributions. First, it formally defines the task of plan recognition in exploratory domains and provides a proof of its NP-completeness. Second, it presents new greedy and complete algorithms for solving the plan recognition problem in these domains, providing a formal complexity analysis of these algorithms and comparing them to existing methods. Third, it is the first work to evaluate plan recognition algorithms on real-world data in the domain of flexible pedagogical software.

We compared two algorithmic approaches for recognizing users' interactions. One of the approaches employed incomplete greedy algorithms to attempt to build plans from the bottom-up. The complexity of one of these algorithms is polynomial in the size of the interaction sequence, while the complexity of the other algorithm is exponential (in the worst case) in the size of this sequence. The second approach converts the recognition process to a Constraint Satisfaction Problem (CSP) using one of two methods. One of these methods builds a complete plan to recognize the entire interaction sequence. The other method works piecemeal in a way that uses subsets of the activity sequence to eliminate infeasible plans before attempting to recognize the entire sequence. This second method was suggested by Quilici et al. (1998) but first tested empirically here. In contrast to the greedy approach, the constraint satisfaction approach is complete, in the sense that if all of the recipes for solving a given TinkerPlots problem exist, and the student solved the problem, the algorithm is guaranteed to find the plan that explains the student's interaction. The complexity of both of the complete methods is exponential in the size of both the interaction sequence and the data set containing ideal solutions.

We conducted a number of empirical studies to evaluate the ability of these algorithms to recognize the plans used to solve TinkerPlots problems. The studies involved two types of settings: adults using TinkerPlots in a lab setting, and middle school students using TinkerPlots in a classroom setting.

The results confirmed that the complete algorithms were able to recognize all plans when the relevant recipes for the TinkerPlots problems existed, and students were able to solve the problems. However, there was a systematic difference between these two empirical settings and their effect on the plan recognition algorithms. For adult data, the complete methods outperformed the heuristic approaches by 25%. For student data, which reflects the intended user population of TikerPlots, this difference was just 4%. In addition, the heuristic approaches were (on average) an order of magnitude faster than the complete approaches for both data sets. We show that the determinant for run-time is the size of the interaction sequence for the heuristic approaches, and the size of the plan data-base for the complete approaches. Lastly, the interaction sequences obtained from middle school students were significantly longer than those of adults, and in general, students' interactions corresponded to complete solutions less often than adults.

These results show that the heuristic algorithms we devised provide a good balance between performance and time in the pedagogical software domain we considered. More generally, they demonstrate the feasibility of using Artificial Intelligence techniques to support the analysis of users' interaction with flexible, open-ended software. Although our study uses one type of software, the techniques presented here are general and can be used to support the analysis of users' interactions for different types of exploratory systems. Our techniques are of value to software designers and researchers who wish to understand the way people learn and use computer software, as well as to teachers.

After describing related work (Section 1.1), we introduce the TinkerPlots software (Section 2), highlighting the properties that characterize an exploratory domain. In Section 3, we describe the formal tools for representing plans in exploratory domains and the abstract problem of recognizing plans relative to idealized recipes for achieving domain goals. We then (Section 4) draw an analogy between this plan recognition problem and grammar recognition, showing that the problem is equivalent to context-free recognition under a variant interpretation of such grammars. The analogy allows a simple proof of the NP-completeness of plan recognition in exploratory domains. We present a variety of plan recognition algorithms for this problem in Section 5, and evaluate their performance empirically on data obtained from users' interactions with TinkerPlots in Section 6, demonstrating the practicality of the best of our algorithms in both coverage and speed.

1.1 Related Work

Plan recognition is a cornerstone problem of AI and a necessary component of many applications such as software help systems (Bauer et al., 1993; Mayfield, 1992), story understanding (Charniak and Goldman, 1993; Wilensky, 1978), and natural language dialogue (Carberry, 1990; Grosz and Sidner, 1990). Early approaches have assumed a goal-oriented agent whose activities were consistent with its knowledge base, and which formed a single encompassing plan (Kautz, 1987; Lochbaum, 1998). A notable exception is Pollack (1987) that allowed for agents to have ill-formed plans about achieving certain goals, and Brown and Burton (1978) that allowed for agents' knowledge to be possibly incorrect. We refer the reader to Carberry (2001) for a detailed account of these approaches and focus this section on more recent works which capture some of the endemic qualities of exploratory domains, namely extraneous actions or mistakes, interleaving of activities, and free order among plan constituents.

We first detail approaches that considered temporal relationships among actions that make up agents' plans. Weida and Litman (1992) proposed a method for recognizing plans that explicitly included ordering constraints in the plan library and suggested various criteria for matching plans to action sequences, assuming that each action is directed at completing one of the plans in the library. Avrahami-Zilberbrand and Kaminka (2005) encoded relationships between action parameters in plans using tree structures and provided methods for plan recognition that traverse the tree in a manner that is temporally consistent with the observations. Another approach to handling temporal relationships in plans derives from the analogy between plan recognition and grammar recognition (Sidner, 1985; Geib and Steedman, 2007). Immediate-Dominance/Linear-Precedence (ID/LP) grammars (Gazdar, 1985) describe languages that allowed for linear precedence and free word ordering over rule constituents. Algorithms for parsing ID/LP grammars, which are analogous to recognizing plans, can at times provide exponential savings as compared to considering every possible order configuration of the rule constituents (Shieber, 1984; Barton, 1985). Pynadath and Wellman (2000) developed a probabilistic grammar for modeling agents' plans that also included their beliefs about the environment. These techniques did not allow for interleaving plans. All reordering among plan constituents in these above works was restricted to local permutation among the constituent actions of sub-plans.

Goldman et al. (1999) proposed a probabilistic model of plan recognition that recognized interleaving actions and output a disjunction of plans — rather than a single hierarchy — to explain an action sequence. It also accounted for missing observations (e.g., not seeing an expected action in a candidate plan makes another candidate plan more likely). The algorithm was generative, that is, with each observation, a pending set of possible hypotheses were generated that were subsequently matched against future observations. Geib and Goldman (2009) have augmented this work to allow to recognize multiple instances of the same plan, in addition to interleaving actions. This work provides a bottom-up algorithm that maintains a distribution over the set of possible explanations matching users’ observations, while not assuming that agent’s top-level goal is known. Our work is distinct from this approach in several ways. First, the settings studied by Geib and Goldman do not account for agents’ extraneous actions, an endemic property of the exploratory domain we consider in this paper. Second, the probabilistic approach used by Geib and Goldman is complete when considering full observation sequences, with a worst-case complexity that is exponential in the size of the grammar. We provide heuristic algorithms that may be exponentially more efficient than complete approaches. Our algorithms are parameter-free, designed for ecologically realistic settings (such as classrooms) in which tuning or learning parameters is difficult because of the effort involved in obtaining large amounts of training data. Third, we show the efficacy of our approach on real-world data obtained from students and adults using pedagogical software, whereas Geib and Goldman use synthetic data.

Several works have used probabilistic reasoning to recognize students’ goals when interacting with pedagogical software. Conati et al. (2002) used Bayesian networks to model students’ interactions with an intelligent tutor, using probabilistic inference to recognize interleaving actions. Albrecht et al. (1998) suggested a probabilistic approach to infer players’ goals as well as their future actions from observation sequences. They used Dynamic Bayesian Networks to compute a posterior distribution over possible goals given players’ actions in the game. They are able to capture agents’ mistakes, but infer the likelihood of a single goal or action, rather than recognizing a hierarchical plan representing the entire action sequence. Quilici et al. (1998) proposed an algorithm for implementing plan recognition as a constraint satisfaction problem but do not evaluate it on real data. We augment this work in several ways. First, by implementing this algorithm on ecologically realistic data, that of adults and middle school students using

pedagogical software. Second, by describing alternative, heuristic approaches to complete algorithms in exploratory domains that provide a balance between completeness and time, and comparing the efficacy of these algorithms to the complete, CSP approach.

Lastly, we will mention work in the intelligent tutoring systems community that has been applied to models of students' learning of mathematics and physics (Corbett et al., 2000; Beck and Woolf, 1998; Anderson et al., 1995; VanLehn et al., 2005). In these domains the tutor is an active participant in the student's learning process and ambiguities or uncertainties about the students' plan of action are resolved by querying the student. By contrast, the TinkerPlots style of educational software allows students to "learn by doing" in an exploratory open-ended manner without explicit guidance by a software tutor. Our approach addresses a different problem, that of non-intrusive recognition of students' activities given their complete interaction histories with the software. Past work on recognition of users' goals with computer systems has focused on fixed, strongly constrained settings such as UNIX command line syntax (Blaylock and Allen, 2005), or applications such as medical diagnosis and email notifications in which users tend to adopt the same goals many times (Bauer, 1996; Horvitz, 1999; Lesh, 1997). In educational domains, goals are constantly evolving to reflect new concepts, and it may be difficult to collect student-specific training data for each type of goal.

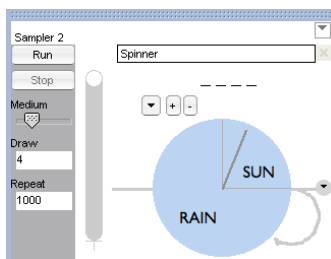
2 The TinkerPlots Domain

TinkerPlots is an educational software system used world-wide to teach students in grades 4 through 8 about statistics and mathematics (Konold and Miller, 2004). It provides students with a toolkit to actively model stochastic events, and to create and investigate a large number of statistical models (Hammerman and Rubin, 2004). As such, it is an extremely flexible application, allowing for data to be modeled, generated, and analyzed in many ways using an open-ended interface.

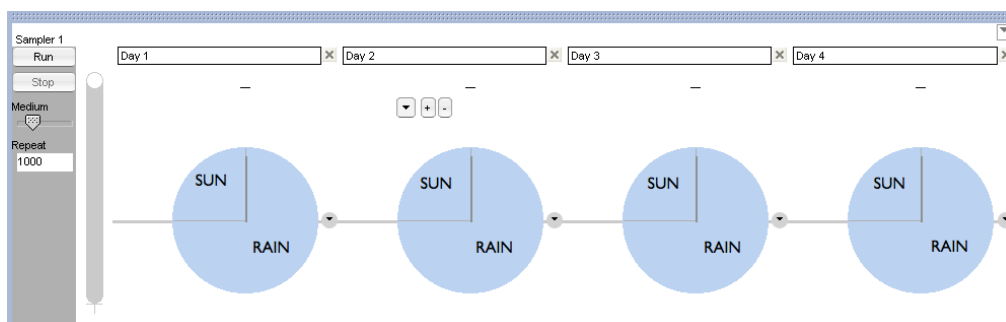
To demonstrate our approach towards recognizing activities in TinkerPlots we will use the following running example, called RAIN.

RAIN: The probability of rain on any given day is 75%. Use TinkerPlots to compute the probability that it will rain on each of the next four consecutive days.

This problem is a simple example drawn from a set of problems posed to students using TinkerPlots in schools and to subjects during our data collection procedure.



(a) Using one spinner

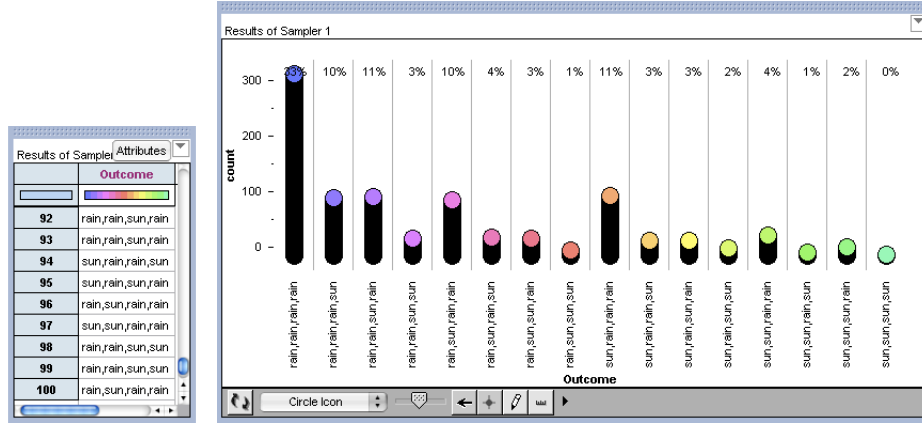


(b) Using four spinners

Figure 1: Two sampler models for RAIN

Two of the several possible approaches towards modeling this problem in TinkerPlots are shown in Figure 1. One uses the same stochastic device multiple times, while the other uses multiple stochastic devices. Figure 1a shows a sampler object containing a single “spinner” device. Devices are added to sampler objects to model distributions. There are several types of devices; spinner devices recall the distribution formed by spinning a dial. The spinner device in the left-hand model contains two possible events, “rain” and “sun”. The likelihood of “rain” is three times that of “sun”, as determined by the surface area of these events within the spinner. Each draw of this sampler will sample the weather for a given day. The number of draws is set to four, making the sampler a stochastic model of the weather on four consecutive days.

Another possible approach to modeling the RAIN distribution is presented in Figure 1b, which shows a sampler with four spinner devices. Each of



(a) Table of Sample Results

(b) Plotted Results

Figure 2: Generating and analyzing data for the RAIN problem.

these devices is a stochastic model of the weather on a given day, and the sampler, set to a single draw, draws once from each device. In both of these approaches, the sampler, along with the contained devices, is a model of the joint probability distribution over the weather for four consecutive days.

When a sampler is run, it generates data that is sampled according to the distribution defined by the parameters of its model. Figure 2a shows a table object holding a portion of the sample generated by either of the sampler models in Figure 1. Each line in the table represents a single repetition of the sampler, consisting of a “sun” or “rain” value for each of four days. Figure 2b shows the end-result of a process in which this data is plotted onto a histogram for the purpose of inferring the likelihood of four consecutive days of rain. There are many other approaches for modeling the RAIN distribution and organizing the resulting data, which we do not show here.

As students interact with TinkerPlots through its engaging direct manipulation interface, they create and modify devices, sample stochastic events, graph the results, modify and retry aspects, in a fluid manner in which all kinds of objects can be manipulated in different orders, and with false starts and retries adding complexity to the exhibited behaviors. The TinkerPlots system is metered to log all the primitive direct manipulation actions of the user.

These logs constitute the trace of the observable behavior of the user.

Our goal is to explain the log in terms of the problem-solving goal that the user was engaged in, such as solving the RAIN problem.

3 Methods and Representation

In this section we introduce representations and algorithms towards describing TinkerPlots activities in a formal way.

3.1 Actions, Recipes and Plans

The nomenclature in this paper follows the foundational planning terminology grounded in philosophy (Bratman et al., 1988; Bratman, 1987; Davidson, 1980). The most fundamental components we define are called *basic actions*, which are atomic, and cannot be decomposed. *Complex actions* describe higher-level, more abstract activities that can be decomposed into sub-actions, which can be basic actions or complex actions themselves. To emphasize the distinction between basic and complex actions, we notate complex actions using an underline notation.

A *recipe* (Pollack, 1990) for a complex action characterizes the sequences of actions that result in successful completion of the action. The recipe for a complex action \underline{C} is a set of sub-actions S and constraints R such that performing the sub-actions under the constraints constitutes completing the complex action. We do not allow recursively defined recipes; i.e, a recipe for the complex action \underline{C} may not hereditarily include \underline{C} in sub-action list S .

The set of *restrictions* R constrains how sub-actions may be completed by expressing relationships over the parameters of the sub-actions that must hold. Restrictions may take the form of any Boolean relation over sub-actions' parameters, which includes mathematical equations and inequalities. A common type of restriction uses inequalities and the **pos** (position) parameter of various actions to limit the order in which these actions must occur. In the absence of ordering restrictions, a recipe is completely free-ordered. Other restrictions may enforce a relationship between object identifiers in several TinkerPlots actions, requiring, for example, that two actions share the same i_s parameter to represent constraints that are imposed on the same sampler object s . To complete complex action \underline{C} according to a recipe, all sub-actions in S must be completed without violating any restrictions in R .

We notate a recipe for complex action \underline{C} with sub-actions $\{s_1, \dots, s_n\}$ and restrictions R as

$$\underline{C} \rightarrow s'_1, \dots, s'_n \text{ where } R$$

where s'_i is the name of the sub-action s_i (with optional subscripts to uniquely identify sub-actions with the same name). The restrictions R use the notation $A[p]$ to refer to the value of a parameter p of some sub-action with name A . (In case multiple sub-actions have the same name, the subscripts are used to disambiguate.) Standard conventions are used to notate multiple inequalities (for example, $a < b < c$).²

A *recipe library* (Bratman et al., 1988) contains the complete set of recipes for all of the complex actions of the domain. Each complex action type may have multiple recipes in the library providing alternatives for its completion.³

A *plan* is a hierarchical construction of basic and complex actions used to complete a complex action called the *root action*. The plan for completing a root action \underline{C} is a tree of parametrized actions rooted at \underline{C} such that each complex action is decomposed into sub-actions according to some recipe in the database.

3.2 Representation of TinkerPlots Activities

The nature of the questions (such as RAIN) for teaching statistical skills in TinkerPlots typically will require students to plan a series of activities to derive answers. Students interact with TinkerPlots through a series of operations that create, modify, or delete objects such as samplers, plots, and tables. Basic actions in TinkerPlots refer to rudimentary operations that can be carried out by a single keystroke or mouse action. It is these instances of basic actions that are logged as the system is used. Examples of basic actions in TinkerPlots include creating a new sampler, generating a random sample, or deleting a plot. Complex actions in TinkerPlots are activities such as adding a spinner with six equally weighted events to a sampler, fitting sampler data to a plot, or solving the RAIN problem. We impute complex

²Our representations of recipes and restrictions are similar to classical planning formalisms such as Hierarchical Task Networks (Ghallab et al., 2004), but do not allow for recursion.

³Other works have used the term “plan library” to refer to complete plan hierarchies from an agent’s root goal down to the basic level actions; see for example Nau et al. (1998). We use the term “recipe library” to refer to a set of recipes.

```

...
AS[is = 11, pos = 5]
ADS[is = 11, id = 2, td = spinner, pos = 6]
ALE[is = 11, id = 2, ie = 1, le = a, pos = 7]
ALE[is = 11, id = 2, ie = 2, le = b, pos = 8]
AS[is = 9, pos = 9]
CEL[is = 11, id = 2, ie = 1, le = rain, pos = 10]
CEL[is = 11, id = 2, ie = 2, le = sun, pos = 11]
CPD[is = 11, id = 2, ss = 1 : 3, pos = 12]
ADS[is = 11, id = 3, td = mixer, pos = 13]
...

```

Figure 3: A snippet of an action sequence taken from a user’s interaction with TinkerPlots

actions to users of the software in our analysis of users’ actions as pursuing plans.

Users’ interactions with TinkerPlots are recorded as finite, chronological sequences of basic actions that are performed by the users. It is these action sequences that constitute the input to the plan recognition algorithm. Figure 3 shows a portion of an action sequence for creating the stochastic component (called a “device”) in the sampler of Figure 1a.⁴ For example, the action `ADS[is = 11, id = 2, td = spinner, pos = 6]` (Add Device to Sampler) refers to the action of adding a device with an identifier `id = 2` and type `td = spinner` to a sampler with an identifier `is = 11`. The `pos` parameter specifies the temporal position of the action within an action sequence; in this case `ADS` is the sixth action performed by the user. Figure 4 provides a key to abbreviations for all actions (in upper-case script) and parameters (in lower-case script) used throughout the paper.

In the TinkerPlots domain, a recipe captures an ideal sequence of actions for performing a particular activity. We represent each basic or complex action type in TinkerPlots as a unique name (such as `ADS` for the basic

⁴This sampler was used to generate the data in Figure 2.

| | |
|--|---|
| <p>Basic Actions</p> <p>ADS = Add Device to Sampler</p> <p>ALE = Add Labelled Event</p> <p>AS = Add Sampler</p> <p>CEL = Change Event Label</p> <p>CPD = Change Probability in Device</p> <p>CSA = Create Sampler with Event A</p> <p>Complex Actions</p> <p><u>AED</u> = Add Event to Device</p> <p><u>CCD</u> = Create Correct Device</p> | <p>Parameters</p> <p>i_d = Device ID</p> <p>i_e = Event ID</p> <p>i_s = Sampler ID</p> <p>l_e = Event Label</p> <p>t_d = Device Type</p> <p>ss = Subsection Size</p> <p>pos = Temporal Position</p> |
|--|---|

Figure 4: Action and parameter abbreviation key

action Add Device to Sampler, or CCD for the complex action Create Correct Device); they are parametrized to describe features of the objects to which an instance of the action refers. We notate an action with its parameters by placing the parameter values, keyed to the parameter names, in brackets after the name. For example, the set of sub-actions for one of the possible recipes for solving RAIN includes creating a sampler that models the weather on four consecutive days, running the sampler, and plotting the results on a graph. Examples of such samplers are shown in Figure 1a and Figure 1b. An example of a graph is shown in Figure 2b.

Figure 5 shows a recipe for the complex action CCD, which creates the sampler shown in Figure 1a. It includes two basic sub-actions, ADS and CPD, as well as two complex AED (Add Event to Device) actions. This recipe also contains several restrictions. The first, an ordering restriction, mandates that a device must be added to a sampler (action ADS) before any events are added to that device (action AED), which in turn must occur before the probability of those events is changed (action CPD). The second and third restrictions require that sampler and device identifiers are consistent across these actions, and the fourth restriction requires that the surface area of events be resized to a 3:1 ratio.

The purpose of giving TinkerPlots problems to students is to test their ability to construct appropriate, applicable models for solving the problem. We treat recipes as idealized descriptions of the use of TinkerPlots to solve problems by constructing a plan for achieving the complex root action for

$$\underline{\text{CCD}} \rightarrow \text{ADS}, \underline{\text{AED}_1}, \underline{\text{AED}_2}, \text{CPD}$$

where

$$\begin{aligned} \text{ADS}[\text{pos}] &\prec \underline{\text{AED}_1}[\text{pos}] \prec \underline{\text{AED}_2}[\text{pos}] \prec \text{CPD}[\text{pos}] \\ \text{ADS}[i_s] &= \underline{\text{AED}_1}[i_s] = \underline{\text{AED}_2}[i_s] = \text{CPD}[i_s] \\ \text{ADS}[i_d] &= \underline{\text{AED}_1}[i_d] = \underline{\text{AED}_2}[i_d] = \text{CPD}[i_d] \\ \text{CPD}[\text{ss}] &= (3 : 1) \end{aligned}$$

Figure 5: A recipe for the CCD (Create Correct Device) complex action

$$\underline{\text{AED}} \rightarrow \text{ALE}, \text{CEL}$$

where

$$\begin{aligned} \text{ALE}[\text{pos}] &\prec \text{CEL}[\text{pos}] \\ \text{ALE}[i_s] &= \text{CEL}[i_s] \\ \text{ALE}[i_d] &= \text{CEL}[i_d] \\ \text{ALE}[i_e] &= \text{CEL}[i_e] \end{aligned}$$

Figure 6: A recipe for the AED (Add Element to Device) complex action

solving the problem. For example, a plan for the complex action CCD is shown in Figure 7. Here, each complex AED action is decomposed into sub-actions ALE and CEL using the recipe shown in Figure 6. The plan was inferred from the student's actions in the sequence shown in Figure 3, and results in the creation of a spinner device, such as the one shown in Figure 1a for solving the RAIN problem.

The flexible nature of TinkerPlots supports exploratory and open-ended use of the software in several ways. First, students may perform extraneous activities that do not play a salient part in the solution to the problem. For example, the action $\text{AS}[i_s = 9, \text{pos} = 9]$ in the student's action sequence of Figure 3 plays no role in the student's plan in Figure 7. Second, students may interleave the sub-actions of different complex actions. For example, in the plan of Figure 7, the $\text{ALE}[i_s = 11, i_d = 2, i_e = 2, l_e = \text{b}, \text{pos} = 8]$ action in position 8 (a sub-action of the complex action AED $[i_s = 11, i_d = 2, i_e = 2, l_e = \text{sun}, \text{pos} = 11]$),

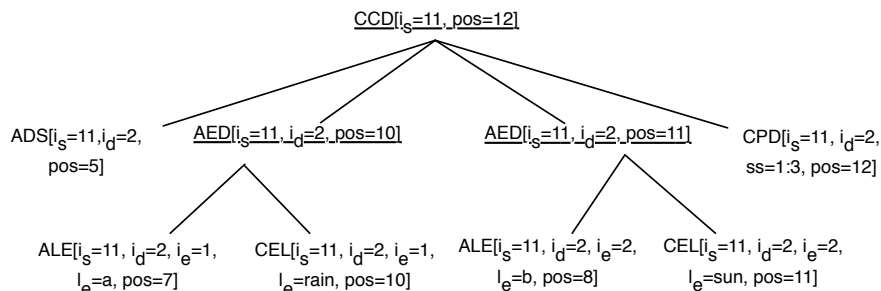


Figure 7: A possible plan for the CCD complex action

temporally occurs among the actions $\text{ALE}[i_s = 11, i_d = 2, i_e = 1, l_e = a, \text{pos} = 7]$ and $\text{CEL}[i_s = 11, i_d = 2, i_e = 1, l_e = \text{rain}, \text{pos} = 10]$ in positions 7 and 10 (sub-actions of the complex action $\text{AED}[i_s = 11, i_d = 2, \text{pos} = 10]$). Lastly, students may make mistakes when solving problems or only succeed in solving part of the problem. The combination of these different properties make it challenging to recognize the plans underlying students' interactions with TinkerPlots, as we argue formally in the next section.

4 Grammars and Complexity

The general problem of recognizing whether a sequence of basic actions embeds a plan that accords with the recipes in the database includes satisfying an arbitrary set of constraints. This problem is at least as complex as constraint satisfaction for the constraint language, which can be NP-hard or worse depending on the particular of the constraint language. But even without the solving of restrictions, the plan recognition problem is NP-complete. Geib and Goldman (2009) show that a related plan recognition problem — involving interleaving and ordering restrictions but not extraneous actions — is NP-hard via a simple reduction from three-dimensional matching to a grammar formalism they call plan tree grammars. The extension to our context, in which extraneous actions in the log are allowed, is straightforward.

In this section, we review and extend the Geib and Goldman proof to derive a complexity result for plan recognition in exploratory domains. For simplicity, we use a simpler formal characterization than the plan tree grammars used by Geib and Goldman. We define a grammar formalism that, like plan tree grammars, allows interleaving, but unlike plan tree grammars,

has no temporal ordering restrictions. This allows us to greatly simplify the description of the formalism. For concreteness, we call the formalism *simple plan grammars*.

A simple plan grammar is structured exactly like a context-free grammar, with a set of terminal and nonterminal symbols including a specified start symbol, plus a set of productions rewriting a nonterminal symbol to a sequence of terminals and nonterminals. As with other grammatical characterizations of planning, the nonterminal symbols correspond to complex actions and the terminal symbols to basic actions. Under this analogy, a recipe corresponds to a grammatical production, a plan to a parse tree, and an action sequence to a string to be parsed. Reconstructing a plan from an action sequence relative to a recipe library would then correspond to parsing a string relative to a grammar (Vilain, 1990).

Although simple plan grammars are structured identically to context-free grammars, the language of a simple plan grammar is defined differently from the corresponding context-free grammar, so as to manifest interleaving and extraneous actions. (Indeed, one can think of a simple plan grammar as an alternate interpretation of the context-free grammar notation.)

The language of a simple plan grammar is defined in two steps. First, we define the *base language* of a simple plan grammar to be the language of the corresponding context-free grammar. The *language* of a simple plan grammar is then the set of all strings containing a subsequence that is a permutation of a string in the base language.

This simple definition captures exactly the reordering, interleaving, and extraneous action aspects of the plan recognition problem, while abstracting away from temporal ordering and other constraints. The reordering and interleaving is captured by the fact that all permutations of the base language strings are in the language of the grammar. The extraneous actions are captured by including supersequences in the language as well, the extra symbols constituting the extraneous actions.

What is not captured by simple plan grammars is the ordering restrictions. It is the ordering restrictions that greatly complicates the definition of plan tree grammars. As we will show, the ordering restrictions are not needed to carry through the NP-hardness proof, and therefore the simpler formalism is sufficient for showing that the plan recognition setting we are considering is NP-hard.

By way of example, consider the following productions:

$$\begin{aligned}
 S &\rightarrow M M M \\
 M &\rightarrow a b c \\
 M &\rightarrow d e f \\
 M &\rightarrow g h i
 \end{aligned}
 \tag{1}$$

If we take these to be the productions of a context-free grammar, the grammar recognizes several strings, including the strings $abcdefghi$ and $abcabcabc$. However, when viewed as a simple plan grammar, it recognizes all supersequences of permutations of these strings, including the strings themselves, but also strings like $adgbehcfi$ or $ihgfedcba$ or $aaaabcdefghiaa$.

It is easy to show that the problem of string recognition for simple plan grammars is NP-complete. We extend the proof of Geib and Goldman (2009), which uses a reduction from the NP-complete problem 3-DIMENSIONAL MATCHING:

3-DIMENSIONAL MATCHING (3DM): Given three identically sized disjoint sets $W = \{w_1, \dots, w_q\}$, $X = \{x_1, \dots, x_q\}$ and $Y = \{y_1, \dots, y_q\}$, and a set $M \subseteq W \times X \times Y$, does there exist a matching consisting of a subset $M' \subseteq M$ of size q such that no two elements of M' agree on any coordinate (that is, all elements of W , X , and Y appear exactly once in M').

This problem was shown to be NP-complete by Karp (1972).

We reduce an instance of 3DM to a simple plan grammar as follows. Given an instance of 3DM, we construct a simple plan grammar with two non-terminals S (the start symbol) and M , and terminal symbols $W \cup X \cup Y$. The productions of the grammar include one for each element $\langle w, x, y \rangle$ of M :

$$M \rightarrow w x y \tag{2}$$

and a production generating q instances of nonterminal M :

$$S \rightarrow \overbrace{M \cdots M}^{q \text{ times}} \tag{3}$$

Note that the base language of the grammar, that is, the language of the productions when viewed as a context-free grammar, comprises strings all of

length $3q$. Thus, if a string in the language includes all of the elements of W , X , and Y , each must occur exactly once.

In addition, we construct the string $s = w_1 \cdots w_q x_1 \cdots x_q y_1 \cdots y_q$ containing each of the elements of W , X , and Y exactly once. We ask whether the string s is admitted by the constructed grammar.

By way of example, the simple plan grammar (1) is exactly the one generated by this construction for the 3DM problem in which $W = \{a, d, g\}$, $X = \{b, e, h\}$, and $Y = \{c, f, i\}$ and where $M = \{\langle a, b, c \rangle, \langle d, e, f \rangle, \langle g, h, i \rangle\}$. The constructed string to recognize would be *adgbehcfi*.

The construction has the property that s is admitted by the constructed simple plan grammar if and only if the corresponding 3DM instance has a solution. (In the example, the constructed string *adgbehcfi* is in the language of the simple plan grammar because its permutation *abcdefghi* is in the base language.) The argument is straightforward, and essentially that of Geib and Goldman (2009), with variation only for the lack of ordering restrictions and the issue of extraneous items.

If there is a solution to the 3DM problem, then there is a subset M' of M that covers all and only the $3q$ elements of $W \cup X \cup Y$. By construction, then, there is a string in the base language that includes all of these elements as well. The constructed string s is a permutation of that string, hence is in the language of the simple plan grammar.

If the string s is in the language of the simple plan grammar, then there is a context-free derivation for some permutation of some subset of the elements of s . Because all strings in the base language are of length $3q$, which is the length of s itself, the base language string must be an improper subset, that is, have exactly the elements of s . But in that case, a solution of the 3DM problem can be read off of the context-free derivation of the base language string. The particular M -productions used in that derivation correspond to the M' subset.

This proof differs from that of Geib and Goldman (2009) in a few ways. First, we do not incorporate ordering constraints in the rule $M \rightarrow wxy$ to require $w \prec x$ and $x \prec y$ as they do. These constraints are not necessary, because by construction the string to be recognized obeys such constraints directly. The same is true of the proof by Geib and Goldman (2009); the ordering constraints are superfluous there too. By observing the superfluity of ordering restrictions for the proof, we allow a simpler grammar setup.

Second, our definition of the language of a simple plan grammar incorporates all supersequences of base language strings, corresponding to allowing

extraneous actions in logs in the plan recognition problem. The original proof was modified to hold even in this context by forcing all base language strings to include exactly $3q$ elements, the same as the constructed string to be recognized, so the issue of supersequences becomes irrelevant. Although the simple plan grammar constructed does admit strings longer than $3q$, they are irrelevant to the argument, as the string to be parsed is of length $3q$. Forcing the string to be of length at least $3q$ is the role of the S production, which has no analog in the Geib and Goldman proof.

Along similar lines, nothing in the constructed grammar enforces the condition that the elements of M chosen are distinct (that is, that no M -production is reused), and no such constraint on the grammar is necessary. If such duplication were to occur, the string generated would have repeated elements as well, but in that case, the derivation will never admit the string to be recognized, which by construction has no repeated elements.

We can conclude, then, that simple plan grammar recognition is NP-hard. The problem is also clearly in NP, as the constructed grammar is polynomial in the size of the 3DM instance, and the context-free derivation for the base language permutation of s serves as a polynomially-sized witness for the recognition problem. Checking that the witness is for a permutation of s is trivially done in polynomial time.

Thus, plan recognition in our model, in which recipes can be interleaved and extraneous actions can be observed, is NP-hard as well. Indeed, this holds whether or not extraneous actions can be observed; they were not made use of in the proof. Similarly, no use was made of recursion in the grammar, so a restriction to non-recursive recipes does not reduce the complexity. Finally, no use of ordering restrictions was made in the proof; satisfying such restrictions makes the recognition problem only more complex.

Given the computational complexity of this plan recognition problem, the question arises as to whether it can be solved in practice for problems of the scope that confront us in real-world cases. We turn to heuristic plan recognition algorithms and their performance in the next section.

5 Plan Recognition Algorithms

In this section we present several plan recognition algorithms that are able to handle the interleaving and extraneous actions that are endemic to exploratory domains such as TinkerPlots. All of these approaches make use of

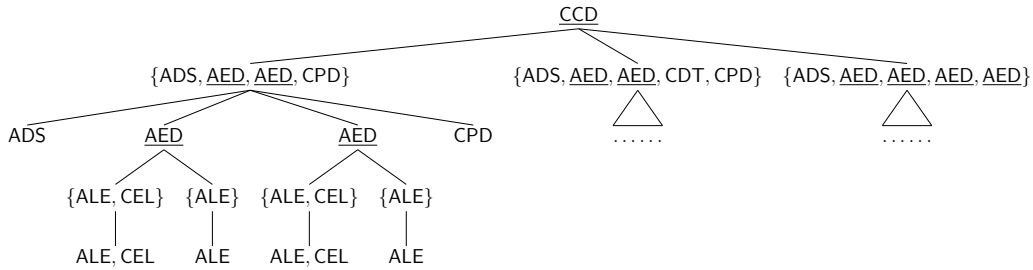


Figure 8: A partial plan tree for the CCD complex action

a structure called a *plan tree* for representing and reasoning about recipes in the database, essentially a search tree for capturing the set of possible plans consistent with the recipe database. A plan tree has two types of nodes: AND nodes, whose children represent actions that must be carried out to complete a recipe, and OR nodes, whose children represent a choice of recipes for completing an action. The root, action C , is an OR node. For each recipe for C , a child AND node is added to the root and labeled with the sub-actions of that recipe. The children of this AND node are the plan trees of each sub-action. A branch terminates when a basic action is reached, as a basic action has no recipe by definition.

A partial plan tree for the CCD action is shown in Figure 8. The AND nodes contain set brackets, while OR nodes do not. Triangles denote unfinished subtrees which were omitted for expository convenience. The plan for creating the spinner object shown in Figure 1a can be found by selecting the leftmost child at each OR node. This resulting plan mirrors the plan shown in Figure 7.

5.1 Greedy Algorithms

We present two greedy algorithms for inferring users' plans. Informally speaking, the algorithms work bottom-up, starting with the user log, and iteratively replacing a set of actions that match the sub-actions of a given recipe by the complex action the recipe implements so as to form a new action list.

A brute-force approach would involve non-deterministically finding all ways in which a complex action might be implemented in the action list. For example, the recipe library for the RAIN problem includes ten recipes and six complex actions. The different recipes for the RAIN problem can form 167,076 possible plans without considering different orderings between

actions.⁵ If we consider all possible orderings within recipes, we get that there are 2,109,182,681,760 possible plans for the RAIN problem. Naively considering each of these possibilities is infeasible.⁶ The heuristic approaches presented in this section make various assumptions about exploratory domains that serve to significantly reduce their complexity as compared to the brute-force method. However, they are incomplete, in the sense that users may construct valid plans that the algorithms fail to infer.

At each step t , the algorithms incrementally build a plan by maintaining an ordered sequence of actions, denoted P_t . The action sequence P_0 , representing the “ground level” of the user’s plan, is denoted as \mathbf{X} . During each step, the algorithms attempt to replace subsets of actions from P_t with the complex actions they represent. Each of the complex actions in P_t is a partial plan that explains some activity in the user’s interaction.

Because our recipe formalization does not allow for recursion, we can define an ordering over all complex actions in the plan library that reflects their depth. Specifically, if $\underline{\mathbf{B}}$ is a constituent sub-action for complex action $\underline{\mathbf{A}}$, then all recipes for $\underline{\mathbf{B}}$ must appear before the recipes for $\underline{\mathbf{A}}$ in the ordering. The heuristic algorithms consider recipes according to this order.

Both greedy algorithms are based on a function `BUILDPLAN` shown in Figure 9 for constructing users’ plans bottom-up. `BUILDPLAN` takes two inputs: An action sequence \mathbf{X} , and a recipe library, \mathbf{R} . The method calls `SORTRECIPES` to topologically sort \mathbf{R} by depth from lowest to highest.

For each recipe R_C for complex action C , the action list P_{t+1} and an open list OL are initialized with the actions in P_t . The algorithm repeatedly tries to find a match for R_c in the open list by calling the function `FINDMATCH(R_C, OL)`, which returns a tuple (M_c, OL) , representing the actions M_C in the match and a modified open list. (The two methods we shall soon present for `FINDMATCH` modify the open list in different ways.) If a match exists, `BUILDPLAN` removes the actions in M_C from P_{t+1} and replaces them with the complex action C , placed where the latest action in the match occurred. (In addition, it marks the actions in M_C in P_t to be the children of action C in P_{t+1} .) Actions are removed from P_t only when a match is

⁵For example, there are 16 possible ways to complete the `CCD` action, because there are two possible recipes in our recipe library for each of the four `MS` actions.

⁶While it is theoretically possible to use string matching to align recipes to the action sequences, a naive approach would need to consider a prohibitive number of possible orderings. The complete approaches (the CSP algorithms) that we describe in the next session essentially perform this matching more efficiently.

```

1: procedure BUILDPLAN(R, X)           ▷ R: a recipe list, X: an action
   sequence
2:    $t \leftarrow 0$ 
3:    $P_0 \leftarrow X$                      ▷  $P_t$ : list of actions at stage  $t$ 
4:   for  $R_C \in \text{SORTRECIPES}(\mathbf{R})$  do   ▷  $R_C$ : a recipe for action  $C$ 
5:      $P_{t+1}, OL \leftarrow P_t$            ▷  $OL$ : an open list
6:      $(M_C, OL) \leftarrow \text{FINDMATCH}(R_C, OL)$    ▷  $M_C$ : a match
7:     while  $M_C$  is not null do
8:       Add  $C$  to  $P_{t+1}$  positioned after last  $a \in M_C$    ▷  $a$ : an action
9:       for all  $a \in M_C$  do
10:        Create a branch from  $C$  in  $P_{t+1}$  to  $a$  in  $P_t$ 
11:        Remove  $a$  from  $P_{t+1}$ 
12:         $(M_C, OL) = \text{FINDMATCH}(R_C, OL, \text{null})$ 
13:    $t \leftarrow t + 1$ 

```

Figure 9: Bottom-up plan recognition method

found, according to the criteria defined in the method FINDMATCH. (We will provide two possible types of criteria in the algorithms below). Once no more matches for R_C can be found, the algorithm moves on to the next recipe, until all the recipes have been considered.

Figure 10 shows several stages of the BUILDPLAN procedure. Edges are shown between complex actions in each step and their constituent actions in the previous step. The first stage, titled P_0 , contains only basic user actions. During the second stage, titled P_1 , two pairs of non-contiguous basic ALE and CEL actions are found to be matches for two complex AED actions. This is an example of interleaving actions, because in P_0 , the ALE constituent of the second AED action is positioned between the ALE and CEL constituents of the first AED action. In the third stage, a match for the CCD action is found, whose sub-actions consist of the first ADS action, both complex AED actions, and the CPD action. Figure 10 defines a structure that is similar to the ideal plan for the CCD action shown in Figure 7, except that it explicitly indicates the interleaving action sequences for the two AED actions, as well as includes redundant actions that were not part of the plan (e.g., AS[is = 9, pos = 8]).

BUILDPLAN is a greedy approach because it does not backtrack. After it chooses a match C for a given recipe R_C , it replaces the actions in P_{t+1} with complex action C without looking ahead to future stages. As a result, the

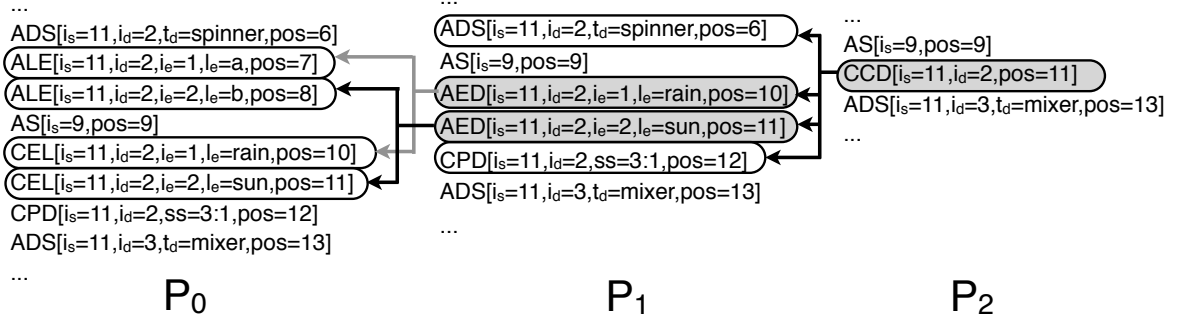


Figure 10: Progression of BuildPlan over three steps

algorithm may fail to find a match for a recipe because a necessary sub-action was committed to another match at an earlier stage. The complexity of the BUILDPLAN algorithm is dominated by the complexity of the FINDMATCH algorithm, call it C_{FM} , discussed in the next section. Let $|\mathbf{R}|$ and $|\mathbf{X}|$ be the number of recipes in \mathbf{R} and the number of actions in the action sequence \mathbf{X} , respectively. Then, BUILDPLAN calls FINDMATCH at most $|\mathbf{X}|$ times per recipe, yielding an overall complexity of $O(|\mathbf{R}| \cdot |\mathbf{X}| \cdot C_{FM})$ for BUILDPLAN.

5.1.1 Matching Algorithms

We present two possible matching algorithms for implementing the FINDMATCH(R_C, OL) process. Both of these make use of the EXTENDS function, a Boolean function that takes as input an action a_P , a partial match M_C , and recipe R_C . It returns true if a_P can be added to M_C , such that (1) a_P corresponds to one of the constituent sub-actions of R_C and is not already in M_C and (2) the addition of a_P to M_C will not violate any of the recipe constraints in R_C . For example, the basic action $\text{ADS}[i_s = 11, i_d = 2, t_d = \text{spinner}, \text{pos} = 6]$ in the action sequence of Figure 3 extends the recipe for CCD shown in Figure 5 given that $M_C = \emptyset$.

The Boolean function FULFILLS(M_C, R_C) returns true if M_C is a complete match for the recipe R_C . We then say that M_C fulfills R_C . Note that M_C

can include both basic and complex actions. For example, the actions

$$\begin{aligned} & \text{ADS}[i_s = 11, i_d = 2, t_d = \text{spinner}, \text{pos} = 6], \\ & \text{AED}[i_s = 11, i_d = 2, t_d = \text{spinner}, \text{pos} = 10], \\ & \underline{\text{AED}[i_s = 11, i_d = 2, t_d = \text{spinner}, \text{pos} = 11]}, \\ & \text{CPD}[i_s = 11, i_d = 2, \text{pos} = 12] \end{aligned}$$

fulfill the recipe for CCD shown in Figure 5.

Both of the matching algorithms choose actions that extend M_C in any order that is allowed by the restrictions of recipe R_C . In particular, the actions in OL may be non-contiguous; this allows the algorithms to capture interleaving plans. However, the two methods differ in the way they update the action list OL as they build a match.

```

1: procedure NOBKTRK( $R_C, OL$ )           ▷  $R_C$ : a recipe,  $OL$ : an open list
2:    $M_C \leftarrow \text{null}$ 
3:   for  $a_P \in OL$  do                       ▷  $a_P$ : an action
4:     if EXTENDS( $a_P, M_C, R_C$ ) then       ▷  $M_C$ : a partial match
5:       Add  $a_P$  to  $M_C$ 
6:       Remove  $a_P$  from  $OL$ 
7:   if FULFILLS( $M_C, R_C$ ) then
8:     return ( $M_C, OL$ )
9:   else if  $M_C$  is null then
10:    return (null,  $OL$ )
11:  else
12:    clear  $M_C$  and goto line 2

```

Figure 11: Algorithm for finding a match without backtracking

The first algorithm, NOBKTRK(R_C, OL), shown in Figure 11, is an extension of an earlier algorithm proposed by Gal et al. (2008). It receives as input R_C , a recipe for some complex action C and an open list OL , which is initially equivalent to the action set in P_{t+1} . NOBKTRK removes actions, one by one, from the open list and places them into a partial match, M_C . Once removed from the open list, these actions will not be reconsidered until a new recipe is provided at step $t + 2$.⁷

⁷We hypothesized that actions that occur late in the interaction process are more salient than actions that occur earlier. However, in practice, traversing the open list in reverse order and increasing order of the pos parameter yielded the same results.

The algorithm is quadratic in the size of the action sequence $|\mathbf{X}|$. To see this, consider that in the worst case, it takes a complete pass over the action list, which is bounded by the size of the action sequence, to fulfill a recipe. Because recipes in TinkerPlots are non-recursive, the number of times a recipe can be fulfilled is also bounded by the size of the action sequence. Therefore, the complexity of NOBKTRK is $O(|\mathbf{X}|^2)$.

```

1: procedure SOMEBKTRK( $R_C, OL$ )           ▷  $R_C$ : a recipe,  $OL$ : open list
2:   return SOMEBKTRKREC( $R_C, OL, \text{null}$ )
3: procedure SOMEBKTRKREC( $R_C, OL, M_C$ )     ▷  $M_C$ : a partial match
4:   if FULFILLS( $M_C, R_C$ ) then
5:     return ( $M_C, OL$ )
6:    $OL' \leftarrow OL$ 
7:   for  $a_P \in OL$  do                       ▷  $a_P$ : an action
8:     remove  $a_P$  from  $OL'$ 
9:     if EXTENDS( $a_P, M_C, R_C$ ) then
10:      Add  $a_P$  to  $M_C$ 
11:      ( $M_C, OL$ ) = FINDMATCH( $R_C, OL', M_C$ )
12:      if FULFILLS( $M_C, R_C$ ) then
13:        return ( $M_C, OL$ )
14:      remove  $a_P$  from  $M_C$ 
15:   return ( $\text{null}, OL$ )

```

Figure 12: Algorithm for finding a match with depth-first search

The second algorithm for finding a match, called SOMEBKTRK, performs a complete depth-first search given a recipe R_C and an open list OL . It defines a sub-function that extends a partial match M_C with a single action, and makes a recursive call to the sub-function. In contrast to the NOBKTRK algorithm, it is complete given a recipe R_C and an open list OL ; that is, it is guaranteed to find a match for R_C if one exists in OL . However, SOMEBKTRK cannot guarantee that a plan is found, because BUILDPLAN itself is greedy. Due to BUILDPLAN's lack of forward-checking or backtracking across time steps, SOMEBKTRK may assign an action to a match during an early step and permanently remove that sub-action from the open list. SOMEBKTRK may later be unable to fulfill a crucial recipe requiring the same sub-action because a match no longer exists in the open list.

$$\begin{aligned}
& \underline{CSA} \rightarrow AS, ADS, ALE \\
& \text{where} \\
& AS[\text{pos}] \prec ADS[\text{pos}] \prec ALE[\text{pos}] \\
& AS[i_s] = ADS[i_s] = ALE[i_s] \\
& ADS[i_d] = ALE[i_d] \\
& ALE[l_e] = 'A'
\end{aligned}$$

Figure 13: A recipe for the CSA (Create Sampler with event “A”) complex action

As an example of the way these two algorithms differ, consider a recipe for a complex action CSA (Create Sampler with Event A) for creating a sampler with one device and a single event labelled “A”, shown in Figure 13. Recall that both NOBKTRK and SOMEBKTRK algorithms extend the current partial match by choosing actions in any order from the interaction sequence that meets the recipe constraints. Given the action sequence shown in Figure 3, the SOMEBKTRK algorithm, which is complete given the recipe for CSA and the action sequence, will find the following match, which includes a sampler with identifier $i_s = 11$ and device with identifier $i_d = 2$.

$$\begin{aligned}
& AS[i_s = 11, \text{pos} = 5], \\
& ADS[i_s = 11, i_d = 2, t_d = \text{spinner}, \text{pos} = 6], \\
& ALE[i_s = 11, i_d = 2, i_e = 1, l_e = 'A', \text{pos} = 7]
\end{aligned}$$

However, the NOBKTRK may decide to add the following actions to the partial match M_C :

$$\begin{aligned}
& ADS[i_s = 11, i_d = 3, t_d = \text{mixer}, \text{pos} = 13], \\
& AS[i_s = 11, \text{pos} = 5]
\end{aligned}$$

The NOBKTRK algorithm will now try to find an $ALE[i_s = 11, i_d = 3, l_e = 'A']$ action which relates to a device with identifier $i_d = 3$, which does not exist in the interaction sequence. Therefore, it will remove the actions in the partial match from consideration. As a result, it will fail to find the $AS[i_s = 11, \text{pos} = 5]$ action in future calls, and will not be able to fulfill the recipe.

To compute the complexity of SOMEBKTRK, let S be the maximum number of sub-actions in any recipe. A first-depth recursive call can be made at

CCD \rightarrow ADS, ALE₁, CEL₁, ALE₂, CEL₂, CPD

where

ADS \prec ALE₁, ALE₂ \prec CPD

CEL₁, CEL₂ \prec CPD

ADS[i_s] = ALE₁[i_s] = CEL₁[i_s] = ALE₂[i_s] = CEL₂[i_s] = CPD[i_s]

ADS[i_d] = ALE₁[i_d] = CEL₁[i_d] = ALE₂[i_d] = CEL₂[i_d] = CPD[i_d]

Figure 14: Expanded version of CCD recipe

most $|\mathbf{X}|$ times. Within each of these recursive calls, at most $|\mathbf{X}| - 1$ actions can remain in OL' . So, at most $|\mathbf{X}| - 1$ second-depth recursive calls can be made for each first-depth recursive call, yielding an overall maximum of $|\mathbf{X}|(|\mathbf{X}| - 1)$ second-depth recursive calls. After $S - 1$ recursive calls of increasing depth have been made, a match must be completed or backtracking must occur. Within each lowest-depth recursive call, there can be at most $|\mathbf{X}| - (S - 1)$ actions left to consider. So, a worst-case complexity of **SOMEBKTRK** is $O(\frac{|\mathbf{X}|!}{S!})$.

5.2 Complete Algorithms

In this section we present two plan recognition algorithms that are complete. Both algorithms work by converting the plan recognition problem into one or more constraint satisfaction problems and using standard techniques for their solution. The conversion makes use of the **EXPAND** function, shown in Figure 15, to convert plans to flat representations containing solely basic actions, called *expanded recipes*. Note that like their conventional counterparts, expanded recipes also include constraints defined over their set of actions. The first complete algorithm performs the conversion naively, while the second use a cascade of conversions to significantly reduce the size of expanded plans.

EXPAND(T_A) takes as input a plan tree T_A for complex action A and returns a set of expanded recipes for A . Each **AND** node represents a possible recipe for its parent node, a complex action. For each **AND** node, **EXPAND**

recursively generates all expanded recipes for each sub-action of the recipe. This algorithm alternates between two sub-procedures, DIRECTSUM and UNION. Given a recipe, the DIRECTSUM procedure computes all possible replacements of complex sub-actions with basic actions. Each time a complex action is replaced, DIRECTSUM ensures that all restrictions involving the complex action are propagated to its sub-actions. For example, consider a single recipe R_A with sub-actions \underline{B} , \underline{C} , and D :

$$\underline{A} \rightarrow \underline{B}, \underline{C}, D$$

with restrictions:

$$\begin{aligned} \underline{B} &\prec D \\ \underline{C} &\prec D \end{aligned}$$

Suppose also that recursive calls of EXPAND have found the expanded recipes for \underline{B} to be $\{E, F\}$ and the expanded recipes for \underline{C} to be $\{G, H\}$ and $\{I, J\}$. In this case, DIRECTSUM will return the following expanded recipes for \underline{A} :

$$\{E, G, H, D\}, \{E, I, J, D\}, \{F, G, H, D\}, \{F, I, J, D\}$$

with each recipe including either the restriction $E \prec D$ or $F \prec D$ in place of the $\underline{B} \prec D$ restriction, and including either the restriction $G \prec D$ and $H \prec D$, or $I \prec D$ and $J \prec D$ in place of the $\underline{C} \prec D$ restriction. Lastly, the UNION sub-procedure takes the union over the expanded recipes generated for each recipe of \underline{A} .

An expanded recipe is a series of basic actions (with associated restrictions) that the user may perform to realize a potential plan. To create an expanded recipe, a path is traversed through the plan tree, beginning at the root and ending with basic actions at the leaves. This path provides a trace of the plan corresponding to the expanded recipe. For example, one expanded recipe can be achieved by traversing the plan tree in Figure 8 and choosing the left-most recipe at each OR node. Notice that the path taken matches the plan in Figure 7. In this expanded recipe, each complex AED action and its restrictions are replaced with two basic actions, ALE and CEL, and corresponding restrictions, as shown in Figure 14.

The complexity of EXPAND is costly in the worst case. Let S be the maximum number of *complex* sub-actions for each recipe, N be the maximum

```

1: procedure EXPAND( $T_C$ ) ▷  $T_C$ : the plan tree for action  $C$ 
2:   ERS[ $C$ ] ←  $\emptyset$  ▷ ERS[ $C$ ]: the expanded recipes for  $C$ 
3:   for all  $r_j$ , a child of  $C$  do ▷  $r_j$ : a recipe
4:     ERS[ $r_j$ ] ←  $\emptyset$ 
5:     for all  $a_i$ , a child of  $r_j$  do ▷  $a_i$ : an action
6:       ERS[ $r_j$ ] ← DIRECTSUM(EXPAND( $T_{a_i}$ ), ERS[ $r_j$ ])
7:       ERS[ $a$ ] ← UNION(ERS[ $a$ ], ERS[ $r_j$ ])
8:     if ERS[ $a$ ] =  $\emptyset$  then
9:       ERS[ $a$ ] ← { $a$ }
10:  return ERS[ $a$ ]

```

Figure 15: Algorithm for generating expanded recipes

number of recipes for a *single* complex action, and C be the number of distinct complex actions. A plan tree has depth of at most $C + 1$, as we do not allow for recursive recipes. At the lowest depth of the plan tree, all actions are basic and do not have recipes. At the second lowest depth, complex actions have at most N expanded recipes, as none of the N recipes contain any complex sub-actions, At the third lowest depth, each recipe for a complex action may contain at most S complex sub-actions, and each sub-action may have at most N recipes. The DIRECTSUM procedure then creates at most N^S expanded recipes per recipe. The UNION procedure collects the expanded recipes resulting from each recipe for that action, resulting in a maximum of $N(N)^S$, or N^{S+1} , recipes. At the fourth lowest depth, each complex action can again have at most N recipes with at most S complex sub-actions in each. Each of these S sub-actions can contain at most $N(N)^S$ expanded recipes. So, the DIRECTSUM and UNION procedures create at most $N(N(N)^S)^S$, or N^{S^2+S+1} , expanded recipes per recipe. Continuing this reasoning, the top-level action can have at most $N^{\sum_{i=0}^{C-1} S^i}$ recipes, yielding an overall complexity of $N^{O(S^C)}$.

5.2.1 Constraint Satisfaction Algorithms

In this section we explain how to combine an expanded recipe and action sequence to create a constraint satisfaction problem (CSP). A solution to the resulting CSP is the plan representing the users' activities. Formally, a CSP is a triple (X, Dom, C) , where $X = \{x_1, \dots, x_n\}$ is a finite set of variables

```

1: procedure CONVERTTOCSP( $\mathbf{E}_A = (S, R), \mathbf{X}$ )    ▷  $\mathbf{E}_A$ : an expanded
   recipe  $S$  and restrictions  $R$  for complex action  $A$ ,  $\mathbf{X}$ : an action sequence
2:   for all  $s \in S$  do                          ▷  $S$ : a set of sub-actions
3:     ADDVARIABLEANDDOMAIN( $s, \mathbf{X}$ )
4:   for all  $r \in R$  do                          ▷  $R$ : a set of restrictions
5:     ADDRESTRICTIONCONSTRAINT( $r$ )
6:   for all  $s \in S$  do
7:     ADDREDUNDANCYCONSTRAINT( $s$ )

```

Figure 16: Converting an expanded recipe and action sequence to a CSP

with respective domains $Dom = \{D_1, \dots, D_n\}$, each a set of possible values for the corresponding variable, $D_i = \{v_1^i, \dots, v_k^i\}$, and a set of constraints $C = \{c_1, \dots, c_m\}$ that limit the values that can be assigned to any set of variables.

The algorithm CONVERTTOCSP, shown in Figure 16, receives as input an expanded recipe E_A and an action sequence \mathbf{X} and returns a CSP. If a solution exists for this CSP, a subset of the actions in \mathbf{X} realize the expanded recipe E_A . We first show how to create variables in the CSP, and we use as a reference Figure 17, which provides a graphical representation of the CSP resulting from the action sequence of Figure 3 and the expanded recipe of Figure 14. We used a graphical layout suggested by Dechter (2003). Note that parameters belonging to actions are not pictured unless they participate in some constraint.

Let $S = \{s_1, \dots, s_n\}$ and R be the set of sub-actions and restrictions in the expanded recipe, respectively. Each action in S becomes a unique variable in the CSP by calling the subroutine ADDVARIABLEANDDOMAIN(s, \mathbf{X}). Based on the expanded recipe, six variables are added at this time: ADS, ALE₁, CEL₁, ALE₂, CEL₂, and CPD. These variables appear, outlined, in the graph of Figure 17. Each variable's domain is then derived from the actions in the action sequence. For each occurrence of action s in the action sequence, a value is added to the domain of s in the CSP. The right-hand box of Figure 17 gives the resulting domain for each variable based on the action sequence.

Lastly, we add restrictions to our CSP. For each restriction r in R over actions (s_1, \dots, s_m) in S , a constraint over the corresponding CSP variables is added to the CSP using the ADDRESTRICTIONCONSTRAINT(r) subroutine. At this point, all restrictions listed in Figure 14 are added, including

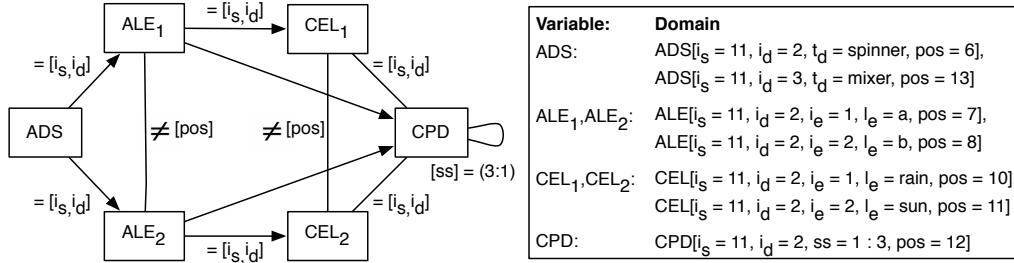


Figure 17: CSP resulting from an action sequence and an expanded recipe

$\text{CPD}[ss] = (3 : 1)$. Directed edges in the figure represent temporal constraints between two variables. Undirected edges represent other parametric constraints. The edge from **ADS** to **ALE₁** expresses the constraint $\text{ADS} \prec \text{ALE}_1$ as well as the constraint $\text{ADS}[i_s, i_d] = \text{ALE}_1[i_s, i_d]$.

For variables corresponding to the same action, additional redundancy constraints are added using the `ADDREDAUNDANCYCONSTRAINT` subroutine. These constraints ensure that such variables are assigned distinct values, as these variables share the same domain. An example is the constraint connecting the **ALE₁** and **ALE₂** variables, which requires that these variable assignments have distinct `pos` parameters.

5.2.2 Brute Force Algorithm

A solution for a CSP provides a match between an expanded recipe and an action sequence. In this section we present two algorithms that use CSPs to output a plan from an action sequenced \mathbf{X} for a desired complex action C given a set of recipes \mathbf{R} .

The first algorithm, shown in Figure 18, takes a brute force approach, calling `EXPAND` to generate each expanded recipe for C , converting it to a CSP and solving the CSP. This algorithm returns the first solution found to the CSP or \emptyset if no solution is found.

The complexity of `CSPBRUTE` can be analyzed in terms of the `FINDMATCH2` and `EXPAND` procedures. Recall that calling `EXPAND` results in at most $N^{O(S^C)}$ expanded recipes, where N is the the maximum number of recipes for a single complex action. In the worst case, all expanded recipes are considered, and for each expanded recipe a CSP solver must be run. The complexity of this CSP solver can be bounded by the complexity of a complete backtracking search, which we have seen to be $\frac{|\mathbf{X}|!}{S!}$. So, an overall

```

1: procedure CSPBRUTE( $T_C, \mathbf{X}$ )  $\triangleright T_C$ : the plan tree for action  $C$ ,  $\mathbf{X}$ : an
   action sequence
2:    $E \leftarrow \text{EXPAND}(T_C)$   $\triangleright E$ : a set of expanded recipes
3:   for all  $e \in E$  do
4:      $\mathbf{C} \leftarrow \text{CONVERTTOCSP}(e, \mathbf{X})$   $\triangleright \mathbf{C}$ : a CSP
5:      $solution \leftarrow \text{SOLVE}(\mathbf{C})$ 
6:     if  $solution \neq \emptyset$  then
7:       return  $solution$ 
8:   return  $\emptyset$ 

```

Figure 18: Brute force algorithm

worst-case complexity of CSPBRUTE is $N^{O(S^C)}O(\frac{|\mathbf{X}|!}{S!})$.

5.2.3 Pruning Algorithm

The second algorithm, shown in Figure 19, takes a more sophisticated approach and traverses the plan tree from the bottom-up. At each OR node, the algorithm determines whether the user completed the corresponding sub-action by either calling CSPBRUTE or referring to the cached result of an earlier CSPBRUTE call. If the user failed to complete that sub-action, the algorithm prunes the relevant recipe, the parent of the current node, from the tree, as the user cannot complete a recipe without completing each sub-action in that recipe. By eliminating branches from the plan tree for the desired complex action C , this pruning process narrows the search space of possible expanded recipes for root action C . This algorithm was suggested by Quilici et al. (1998).

The CSPPRUNE method calls the CSPBRUTE algorithm once per distinct complex sub-action. Let C again represent the number of distinct complex sub-actions in our recipe list. Then, the worst-case complexity of the pruning method is $N^{O(S^C)}O(\frac{C|\mathbf{X}|!}{S!})$. The worst case occurs when the user has completed each of the C complex actions, causing no potential expanded recipes to be eliminated. However, we hypothesized that users would be likely to solve TinkerPlots problems just once, and therefore some complex actions within the plan trees would not be matched in the action sequence.

```

1: procedure CSPPRUNE( $C, \mathbf{R}, \mathbf{X}$ )  ▷  $\mathbf{R}$ : a recipe library,  $\mathbf{X}$ : an action
   sequence
2:    $T_C \leftarrow$  CREATERECIPE TREE( $C, \mathbf{R}$ )  ▷  $T_C$ : the plan tree for action  $C$ 
3:   Perform a bottom-up traversal of  $T_C$ .
4:   for each OR node representing action  $A$  do
5:     if  $A$  has not been cached then
6:       Cache  $solution[A] \leftarrow$  CSPBRUTE( $T_A, \mathbf{X}$ )
7:     if  $solution[A] = \emptyset$  then
8:       Prune parent of  $A$  from  $T_A$ 
9:   return  $solution[C]$ 

```

Figure 19: Bottom-up algorithm

6 Evaluation

The plan recognition problem inherent in domains such as TinkerPlots, where agents are engaged in exploratory behavior with false starts in addition to successful plan construction, leads, as we have shown, to an NP-complete computation. The algorithms presented in the previous section, both incomplete and complete, were intended to allow solution of real-world problems in practice, despite the complexity issue inherent in the problem. To determine their real-world performance, we collected actual logs of TinkerPlots usage on standard pedagogical problems, and compared the algorithms’ coverage and performance. The results show that the best of the algorithms has excellent coverage and practical performance.

6.1 Experimental Design

We collected interaction sequences of people’s interaction with TinkerPlots in two different settings. The first setting included 12 adults with a wide variety of educational backgrounds, ranging from some high school to some post-graduate education. The second setting included 12 eighth grade students in a middle school in Cambridge MA.⁸ Each adult subject received an identical 30-minute tutorial about TinkerPlots and was then asked to complete four problems in succession; these problems are detailed in Appendix A. Students

⁸Appropriate IRB approval was obtained for both settings, and parental consent was obtained for the data collected from the eighth graders.

were given a slightly longer 45 minute demonstration of the software and were asked to solve two of the four problems. User logs and videos of the users' screens were recorded for all user sessions. TinkerPlots is equipped with a logging facility that records the basic actions that make up users' action sequences. To evaluate the various plan recognition algorithms, we manually traced the videos of their interaction with TinkerPlots. We noted whether each problem was solved, and we constructed the (possibly multiple) plans used to solve the problem. We define a recognition algorithm to be "correct" if the plans that it outputs exactly corresponds to the plans constructed from the videos, or if it fails to output a plan when the student did not successfully complete the problem as determined by an expert.⁹ If a user has solved a problem in several different ways, a recognition algorithm is deemed correct if it recognizes any of these solutions.

We created a set of recipes to the TinkerPlots problems in our study to serve as input to the plan recognition algorithms. Our purpose in this study was to evaluate algorithms for matching these ideal solutions with the appropriate basic level actions in users' interaction sequences. Therefore, the recipes we manually constructed were created prior to the collection of and were not informed by the data of people's interactions with TinkerPlots. Rather, they represented what we perceived a priori to be a broad range of possible solutions for TinkerPlots problems. Ultimately, our database contained recipes sufficient to explain all but three user interactions. The lack of inclusion of these recipes is discussed in Section 6.2.2. The recipe library used to run the recognition algorithms on each problem consisted of those recipes that were constructed for the problem. This corresponds to knowing which TinkerPlots problem students are trying to solve. This assumption is logical in the context of pedagogical software.¹⁰

6.2 Results and Discussion

We compared the performance of the four recognition algorithms presented, called the NOBKTRK, SOMEBKTRK, CSPBRUTE, and CSPPRUNE tech-

⁹The domain expert was a researcher of educational technology who has worked with TinkerPlots for several years. For each action sequence, the expert was shown a movie of the desktop of the user, as well as the plan outputted by the algorithm.

¹⁰Dropping this assumption corresponds to running our algorithms using a recipe data base that contains the complete set of recipes for all problems. This may affect correctness for the greedy algorithms, but not for the complete algorithms.

| | # possible plans | S | N | C |
|-----------|------------------|---|---|----|
| ROSA | $O(2^6)$ | 4 | 2 | 13 |
| RAIN | $O(2^{17})$ | 4 | 2 | 6 |
| EARRINGS | $O(2^{11})$ | 4 | 2 | 6 |
| SEATBELTS | $O(2^{12})$ | 4 | 4 | 6 |

Table 1: Number of possible plans per problem, maximal number of complex actions per recipe (S), maximal number of recipes for a single complex action (N), maximal number of distinct complex actions (C).

| Data | Problem | Algorithm | | | |
|----------------------|---------------------|-----------|-----------|-----------|-----------|
| | | NOBKTRK | SOMEBKTRK | CSPBRUTE | CSPPRUNE |
| Adults + Students | ROSA (23) | 78% (18) | 87% (20) | 100% (23) | 100% (23) |
| | RAIN (16) | 88% (14) | 88% (14) | 100% (16) | 100% (16) |
| Adults | SEAT BELTS (12) | 42% (5) | 67% (8) | 75% (9) | 75% (9) |
| | EARRINGS (11) | 91% (10) | 100% (11) | 100% (11) | 100% (11) |
| | Overall (62) | 76% (47) | 85% (53) | 95% (59) | 95% (59) |

Table 2: Accuracy of recognition algorithms by percentage. Parenthesized numbers are the number of logs.

niques.¹¹ We analyzed the user logs corresponding to the problems outlined in Appendix A. Table 1 lists features for each problem that affect the complexity analysis of Section 5.1.1.

The analyzed user logs ranged in length from 14 to 80 actions. The average length of an interaction sequence for problems collected from adult subjects was 35 actions. Adults solved the assigned problems 70% of the time. In contrast, the average length of an interaction sequence for problems collected from students was 68 actions. Students solved the assigned problems 60% of the time. Also, people engaged in exploratory behavior using the software. For example, there were on average 15 exogenous actions in each problem that was obtained from adults.

Table 2 shows the accuracy of the recognition algorithms on data collected for students and adults. As shown by the Table, the heuristic algorithms NOBKTRK and SOMEBKTRK were correct for 47 of 62 (76%) and

¹¹We used the python-constraint package created by Gustavo Niemeyer and available at <http://labix.org/python-constraint> to implement our constraint satisfaction algorithms.

53 of 62 (85%) interactions, respectively. The heuristic algorithms were outperformed by both CSPBRUTEFORCE and CSPPRUNE algorithms, which performed correctly for 59 of 62 (95%) interactions. The incorrect inferences in all algorithms were false negatives, that is, the algorithms were unable to find solutions existing within the interaction sequence. All of the solutions outputted by all of the algorithms matched the expert’s opinion regarding the activities used by the students.

The constraint satisfaction algorithms are guaranteed to find users’ plans if they exist, provided that the relevant recipes are contained in the recipe database. Each of the three incorrect inferences of these algorithms can be traced to recipes missing from the database. In addition to incorrect inferences due to missing recipes, the incomplete approaches suffered from prematurely committing to a match without being able to backtrack. An example of the latter case occurred frequently within the SEAT BELTS problem (see Appendix). This problem required the user to construct a conditional distribution representing the fact that people wearing seat belts are less likely to be hurt in an accident. Some users created a sampler with the wrong distribution parameters and proceeded to correct these parameter settings. The incomplete approaches tended to match sub-recipes for creating a sampler with actions corresponding to the wrong parameters. Consequently, they failed to find matches for the sampler construction recipe. In contrast, the CSP algorithms were able to backtrack and pick the right match for the expanded sampler construction recipe.

6.2.1 Performance of Algorithms on Ecologically Realistic Data

There was a significant difference in the performance of the heuristic algorithm on data that was obtained from adults and the data that was obtained from middle school students. Figure 20 shows the average accuracy of the plan recognition algorithms on student data (left, 19 instances) and adult data (right, 20 instances) from two problems, ROSA and RAIN.¹² The Figure details the performance of these algorithms when correctly identifying successful solutions (true positives), correctly identifying failed solutions (true negatives), and incorrectly identifying failed solutions (false negatives).

As shown in Figure 20, both of the heuristic algorithms were always able to recognize unsolved problems (true negatives). There were 6 such cases

¹²We performed analysis on these two problems for which we have both student and adult data.

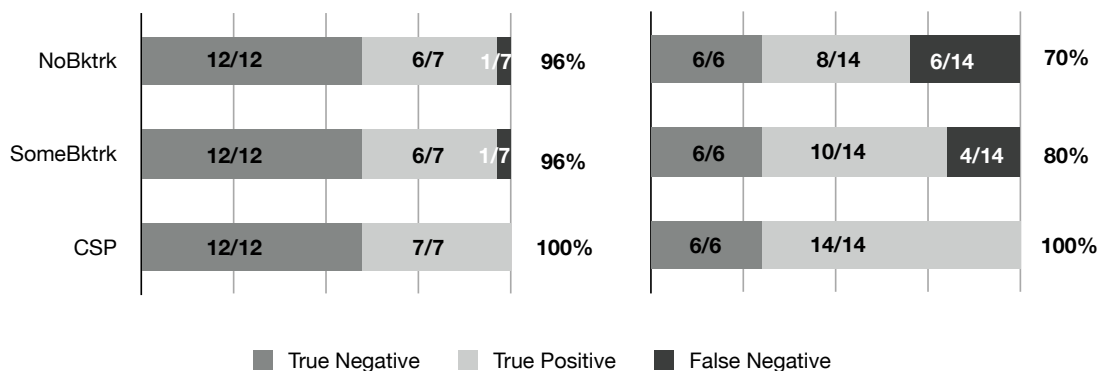


Figure 20: Performance for data obtained from students (left) and adults (right)

for adult data and 12 cases for student data). More generally, for student data the accuracy of both heuristic algorithms averaged 96%, just 4% lower than the accuracy of the complete CSP algorithms for this data. However, for adult data, the average accuracy of the heuristic algorithms was 75%, significantly lower than that the average accuracy of the CSP algorithms. This is also apparent from Table 3, which compares the precision and recall measures of the different approaches on student and adult data. All of the algorithms achieved perfect precision, because there were no false positive classifications, as we described above. The complete CSP algorithms both achieved perfect recall because they did not incur false negatives. For the heuristic approaches, both NOBKTRK and SOMEKTRK algorithms achieved lower recall measures than the CSP approaches for both student and adult data. However, the average recall measure achieved by the heuristic approaches for student data was significantly higher than the average recall they achieved for adult data.

To explain the effect of the type of the setting on the various plan recognition algorithms, we compared the number of times in which students and adults were able to solve the problems that were assigned to them. As shown in Figure 20 students were able to solve TinkerPlots problems less often than adults (for students, 12 data instances represented failures out of a total of 19 instances; for adults, 6 data instances represented failures out of a total of 20 instances). We attribute this to the fact that students were more likely to engage in exploratory behavior or to make mistakes while using TinkerPlots,

| Data | Algorithm | | |
|---------|-----------|-----------|--------|
| | NOBKTRK | SOMEBKTRK | CSP |
| Student | (1, 0.86) | (1, 0.86) | (1, 1) |
| Adult | (1, 0.57) | (1, 0.71) | (1,1) |

Table 3: Precision (left) and recall (right) measures for student and adult data

as is attested by their longer interaction sequences. Thus, one explanation for the success of the heuristic algorithms when analyzing student data is that there were more true negative classifications for student data, and these were always recognizable by our recognition algorithms.

However, the performance of the heuristic algorithms on student data cannot be attributed solely to the fact that there were more true negative classifications for student data. Both of the heuristic algorithms also achieved higher accuracy rates for students than for adults on data instances representing true positives. As shown by Figure 20, for student data, the heuristic algorithms were able to recognize 6 out of 7 true positive instances (average accuracy 85%), but were only able to recognize 8 and 10 out of 14 true positive instances (average accuracy 65%) for adult data. Thus, both NOBKTRK and SOMEBKTRK algorithms were better at recognizing successful and unsuccessful plans for students than for adults. These results provide additional support for the applicability of these algorithms, showing that they are particularly suitable for data that was obtained in an ecologically realistic fashion, rather than a lab setting.

6.2.2 Limitations of Approaches

A significant hurdle to accurate recognition is the difficulty of expressing certain types of user strategies as recipes. As previously mentioned, both the greedy and constraint satisfaction algorithms were susceptible to the fact that the recipe database failed to capture all of the possible ways in which people solved problems. Because we do not allow recipes to be defined recursively, our recipes can only specify that a fixed, rather than variable, number of actions occur. In the RAIN problem, for instance, there are several possible samplers that model the probability of rain as 75%. A natural solution is one (as in Figure 1a) in which there are unique events for “rain” and “sun”,

causing “rain” to be weighted as three times more likely than “sun”. In actuality, any proportional number of “rain” and “sun” events would suffice, as long as “rain” is three times more likely than “sun”. This limitation prevented our library from succinctly expressing the infinite number of permissible strategies – three of which were used by students – for the SEAT BELTS problem.

Although there were no false-positive classifications in our experimental sessions, they can happen in theory. One reason for this is that the recipe language cannot express actions that must not occur. For example, the ROSA problem requires the user to create a device with four events labelled *R*, *O*, *S*, and *A*. If a fifth event were added, the sampler would no longer be correct. However, all of our recognition algorithms would incorrectly match the actions corresponding to the four previous values with the recipe for creating a sampler.

It is possible to construct plans using recipes that cannot happen in practice because they are disallowed by the software. For example, two interleaved actions may add and delete each other’s preconditions.

Finally, we note that our recipes are not designed to describe approximately correct user approaches. For example, one problem required users to add two identically labeled events to a device, without specifying what that label must be. One user failed to complete this task according to the recipe database because he or she used event labels “pierced” and (sic) “piecred”. Despite this oversight, a teacher would likely consider this strategy a successful one. However, our algorithms were unable to distinguish this mistake from other, more conceptual mistakes. One way to overcome such difficulties is to search for the “closest” plan for explaining users’ TinkerPlots activities in terms of recipes, rather than searching for only complete and correct plans. This can be done by using flexible CSP solvers, which search for solutions that minimize the number of violated constraints.

6.2.3 Performance Considerations

In this section we compare the performance incurred in practice by the four recognition algorithms we have presented (measured as run-time on a commodity computer). By working to recognize interleaving plans in user logs containing incorrect strategies and up to 80 actions, we test whether the worst-case exponential complexity presents a significant barrier to real-world plan recognition. In Table 4 we present run-times for each of our four algo-

| Problem | Algorithm | | | |
|----------------|------------------|-----------|----------|----------|
| | NOBKTRK | SOMEBKTRK | CSPBRUTE | CSPPRUNE |
| ROSA | 0.15 | 3.34 | 0.11 | 0.09 |
| RAIN | 0.35 | 0.54 | 131.54 | 21.88 |
| SEAT BELTS | 0.02 | 0.34 | 6.07 | 8.42 |
| EARRINGS | 0.01 | 1.82 | 6.92 | 3.99 |
| Overall | 0.13 | 1.54 | 36.75 | 8.02 |

Table 4: Average runtime (in seconds) of recognition algorithms

rithms organized by increasing theoretical computational complexity. These results are averaged over all action sequences for each problem.

As shown in the table, the average run-time of the heuristic approaches, which employed limited or no backtracking, was significantly faster than that of the complete approaches. Within the complete approaches, the pruning algorithm was significantly faster than the brute force algorithm. This is because there were “dead-ends” in the recipe tree that the pruning algorithm was able to exploit. As can be seen in the table, CSPPRUNE outperformed the CSPBRUTEFORCE algorithm, and in reasonable time, despite having greater worst-case complexity. However, in cases where the pruning approach fails to eliminate branches from the plan-trees, it may turn out to be slower than the brute force approach. For example, as shown in the Table, for the SEAT BELTS problem, the CSPBRUTEFORCE algorithm was faster than the CSPPRUNE algorithms.

To examine the relationship between runtime and student interaction length, we measured the correlation between these variables for each algorithm. The heuristic algorithms NOBKTRK and SOMEBKTRK showed positive correlations between runtime and interaction length of .752 and .508, respectively. The CSPBRUTEFORCE algorithm showed a negative correlation of -.333, while CSPPRUNE showed a very weak, positive correlation of .050. Also, as shown in Table 4, the complete approaches were significantly slower when running on RAIN than when running on the other problem. The number of possible plans for RAIN, as shown in Table 1, was significantly higher than the number of possible plans for the other problems.

These results support our complexity findings that the bottleneck of the incomplete algorithms is the size of the user’s interaction log, while the bottle-

| Problem | CSP _{BRUTE} | | CSP _{PRUNE} | |
|----------------|----------------------|------------------|----------------------|------------------|
| | CSPs | Variables | CSPs | Variables |
| ROSA | 19 | 18 | 14 | 4 |
| RAIN | 9300 | 29 | 141 | 9 |
| SEAT BELTS | 347 | 23 | 298 | 11 |
| EARRINGS | 292 | 21 | 176 | 14 |
| Overall | 2490 | 23 | 157 | 10 |

Table 5: Average number of CSPs and variables for constraint satisfaction algorithms

neck of the complete algorithms is the complexity of the recipe library. These results also demonstrate the applicability of our algorithms to be used in actual classrooms. As shown in Table 1, the number of possible plans for some TinkerPlots problems is is very large. Despite this fact, the almost-perfect record of the heuristic approach, and the short runtime of the complete approaches on these instances speaks well for their overall performance.

Lastly, to further compare the performance of our constraint satisfaction algorithms, Table 5 presents two additional statistics: the average number of CSPs built per log and the average number of variables contained in each CSP. CSP_{PRUNE} outperforms CSP_{BRUTEFORCE} in each category, building fewer and smaller CSPs, on average.

7 Future Work and Conclusion

This paper investigated a class of plan recognition problems for domains in which agents engage widely in exploratory behavior. We showed that constraint satisfaction algorithms are a viable and practical approach for plan recognition in one such domain, that of an educational software application. These algorithms were able to correctly capture users’ plans in real-world logs of users’ sessions in reasonable time despite the theoretical worst-case behavior and the flexible nature of the software. The algorithms compared favorably to faster but greedier approaches.

This work is a first step towards a pedagogical agent that is truly collaborative, in the sense that it provides useful machine-generated support

to teachers and students. For teachers, this support consists of information about students' performance both during and after class. For example, presenting teachers with a visualization of students' plans will convey whether and how students solved a particular problem more quickly than would be possible if they had to analyze snapshots. Teachers can also benefit from the fact that our algorithms capture false starts and incorrect solutions, alerting them to mistakes and misconceptions by the students. Existing systems for assessing students' performance with pedagogical software work in highly constrained settings, and report simple statistics, such as the number of correct answers solved (Feng et al., 2009). Our work extends such systems to exploratory domains, in which students' performance can be explained in part by inferring their plans. In a recent user study we conducted with teachers using TinkerPlots in the classroom, teachers favored the plan-based presentation to other types of visualizations, such as seeing selected snapshots of the students' work.

We are currently extending these results in two ways. First, we are developing methods for presenting plan recognition output to teachers in order to provide them with a broad and organized view of students' activities. Second, we are evaluating the ability of the algorithms to generalize to a different pedagogical software system for teaching chemistry to college students (Yaron et al., 2009).

This work raises several new opportunities for involving the use of plan recognition algorithms as a basis for building intelligent tutors that will augment existing software tools for mathematics education. These collaborative tutors will provide machine-generated support that decides when and how to intervene with students based on teachers' feedback to the plans inferred by the system. They will contribute to the thoughtful analysis of probabilistic models by students and increased ability of teachers to identify those students who would benefit from teacher advice. Together these abilities should lead to improvements in both teaching and learning.

8 Acknowledgments

This work was supported by NSF grant number REC-0632544. We are indebted to Craig Miller for developing the TinkerPlots logging facility. Thanks very much to the anonymous reviewers for their helpful comments. Thanks to Elif Yamangil for assisting with the development of the greedy algorithm,

and to Cliff Konold for helpful discussions. A special thanks goes to Ofra Amir for reading and commenting on previous drafts of this work.

A Experimental Problems

We detail the four TinkerPlots questions posed to subjects and considered in our empirical evaluation of the algorithms.

ROSA: Jessica has 4 letters printed on cards: R, O, S, and A. After mixing them up, she blindly picks the 4 letters one at a time and arranges them in line in the order she chose them. Build a TinkerPlots model and use it to help you estimate the probability of Jessica spelling the word ROSA.

RAIN: There is a 75% chance of rain for each of the next 4 days. Build a TinkerPlots model and use it to help you estimate the probability of getting rain on all 4 days.

SEAT-BELTS: If you get into an accident, you are much less likely to be injured if you are wearing your seat belt. Build a TinkerPlots model for people that

1. are either wearing seat belts or not,
2. then are either injured in accident or not injured in an accident.

Design your factory so that the people wearing seat belts are less likely to get injured than those not wearing them. In your model, what is the probability of people being injured in an accident? What is the probability of being injured in an accident when you are wearing a seat-belt?

EARRINGS: Build a TinkerPlots factory that

1. makes people that are girls or boys,
2. then either pierces their ears or not.

According to your model, what is the probability that

1. a boy has a pierced ear?

2. a girl has a pierced ear?

According to your model, approximately what fraction of people you meet on the street will have pierced ears?

References

- D.W. Albrecht, I. Zukerman, and A.E. Nicholson. Bayesian models for key-hole plan recognition in an adventure game. *User modeling and user-adapted interaction*, 8(1):5–47, 1998.
- J. R. Anderson, A. T. Corbett, K. R. Koedinger, and R. Pelletier. Cognitive tutors: Lessons learned. *The Journal of Learning Sciences*, 4(2):167–207, 1995.
- D. Avrahami-Zilberbrand and G.A. Kaminka. Fast and complete symbolic plan recognition. In *Proceedings of the International Joint Conference on Artificial Intelligence*, volume 14, 2005.
- T. Babaiian, B. J. Grosz, and S. M. Shieber. A writer’s collaborative assistant. In *Intelligent User Interfaces Conference*, pages 7–14, 2002.
- G. E. Barton. On the complexity of ID/LP parsing. *Computational Linguistics*, 11(4):205–218, 1985. ISSN 0891-2017.
- M. Bauer. Acquisition of user preferences for plan recognition. In *Proceedings of the Fifth International Conference on User Modeling*, pages 105–112, 1996.
- M. Bauer, S. Biundo, D. Dengler, J. Koehler, and G. Paul. A logic-based tool for intelligent help systems. In *Proc. 13th International Joint Conference on Artificial Intelligence (IJCAI)*, 1993.
- J. E. Beck and B. P. Woolf. Using a learning agent with a student model. In *Proc. of 4th international conference on Intelligent Tutors*, 1998.
- N. Blaylock and J. Allen. Recognizing instantiated goals using statistical methods. In *Workshop on Modeling Others from Observations*, pages 79–86, 2005.

- M. E. Bratman. *Intention, Plans, and Practical Reason*. Harvard University Press, Cambridge, MA, 1987.
- M. E. Bratman, D. J. Israel, and M. E. Pollack. Plans and resource-bounded practical reasoning. *Computational intelligence*, 4(3):349–355, 1988.
- J.S. Brown and R.R. Burton. Diagnostic models for procedural bugs in basic mathematical skills. *Cognitive science*, 2(2):155–192, 1978.
- H.H. Bui. A general model for online probabilistic plan recognition. In *Proc. 18th International Joint Conference on Artificial Intelligence (IJCAI)*, 2003.
- S. Carberry. *Plan Recognition in Natural Language Dialogue*. MIT Press, 1990.
- S. Carberry. Techniques for plan recognition. *User Modeling and User-Adapted Interaction*, 11(1):31–48, 2001.
- E. Charniak and R. P. Goldman. A Bayesian model of plan recognition. *Artificial Intelligence*, 64(1):53–79, 1993.
- P.R. Cohen, C.R. Perrault, and J.F. Allen. Beyond question-answering. In W. Lehnert and M. Ringle, editors, *Strategies for Natural Language Processing*, pages 245–274. 1981.
- C. Conati, A. Gertner, and K. VanLehn. Using Bayesian networks to manage uncertainty in student modeling. *Journal of User Modeling and User-Adapted Interaction*, 12(4):371–417, 2002.
- A. Corbett, M. McLaughlin, and K. C. Scarpinato. Modeling student knowledge: Cognitive tutors in high school and college. *User Modeling and User-Adapted Interaction*, 10:81–108, 2000.
- D. Davidson. Intending. In *Essays on Actions and Events*, pages 83–102. Clarendon Press, 1980.
- R. Dechter. *Constraint Processing*. Morgan Kaufmann, 2003.
- M. Feng, N. Heffernan, and K. Koedinger. Addressing the assessment challenge with an online system that tutors as it assesses. *User Modeling and User-Adapted Interaction*, 19(3):243–266, 2009. ISSN 0924-1868.

- Y. Gal, E. Yamangil, A. Rubin, S. M. Shieber, and B. J. Grosz. Towards collaborative intelligent tutors: Automated recognition of users' strategies. In *Proceedings of Ninth International Conference on Intelligent Tutoring Systems (ITS), Montreal, Quebec, 2008*.
- G. Gazdar. *Generalized phrase structure grammar*. Harvard Univ Pr, 1985.
- C. W. Geib and M. Steedman. On natural language processing and plan recognition. In *Proceedings of the 20th international joint conference on Artificial intelligence (IJCAI)*, pages 1612–1617, San Francisco, CA, USA, 2007. Morgan Kaufmann Publishers Inc.
- C.W. Geib and R.P. Goldman. A probabilistic plan recognition algorithm based on plan tree grammars. *Artificial Intelligence*, 173(11):1101–1132, 2009.
- M. Ghallab, D.S. Nau, and P. Traverso. *Automated Planning: theory and practice*. Morgan Kaufmann Publishers, 2004.
- R.P. Goldman, C.W. Geib, and C.A. Miller. A new model of plan recognition. In *Proc. 15th Conference on Uncertainty in Artificial Intelligence (UAI)*, 1999.
- B. J. Grosz and C. L. Sidner. Plans for discourse. *Intentions in Communication*, pages 417–444, 1990.
- J. K. Hammerman and A. Rubin. Strategies for managing statistical complexity with new software tools. *Statistics Education Research Journal*, 3(2):17–41, 2004.
- E. Horvitz. Principles of mixed-initiative user interfaces. In *Proc. of ACM SIGCHI Conference on Human Factors in Computing Systems*, 1999.
- R. M. Karp. Reducibility among combinatorial problems. In R. E. Miller and J. W. Thatcher, editors, *Complexity of Computer Computations*, pages 85–103. New York, 1972.
- H. A. Kautz. *A formal theory of plan recognition*. PhD thesis, University of Rochester, 1987.
- C. Konold and C. Miller. *TinkerPlots Dynamic Data Exploration 1.0*. Key Curriculum Press, 2004. URL <http://www.keypress.com/x5715.xml>.

- N. Lesh. Adaptive goal recognition. In *Proceedings of the 15th International Joint Conference on Artificial Intelligence*, pages 1208–1214, 1997.
- K. E. Lochbaum. A collaborative planning model of intentional structure. *Computational Linguistics*, 4(525–572), 1998.
- J. Mayfield. Controlling inference in plan recognition. *User Modeling and User-Adapted Interaction*, 2(1):55–82, 1992.
- D.S. Nau, S.J.J. Smith, K. Erol, et al. Control strategies in HTN planning: Theory versus practice. In *Proceedings of the National Conference on Artificial Intelligence*, pages 1127–1133, 1998.
- M.E. Pollack. Some requirements for a model of the plan inference process in conversation. *Communication failure in dialogue and discourse: detection and repair processes*, 1987.
- M.E. Pollack. Plans as complex mental attitudes. *Intentions in communication*, 1990.
- D.V. Pynadath and M.P. Wellman. Probabilistic state-dependent grammars for plan recognition. In *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence*, pages 507–514, 2000.
- A. Quilici, Q. Yang, and S. Woods. Applying plan recognition algorithms to program understanding. *Automated Software Engineering*, 5(3):347–372, 1998.
- S. Reddy, Y. Gal, and S. M. Shieber. Recognition of users’ activities using constraint satisfaction. In *Proceedings of the First and Seventeenth International Conference on User Modeling, Adaptation and Personalization*, 2009.
- K. Ryall, J. Marks, and S. M. Shieber. An interactive constraint-based system for drawing graphs. In *Proceedings of the 10th Annual Symposium on User Interface Software and Technology (UIST)*, 1997.
- S.M. Shieber. Direct parsing of ID/LP grammars. *Linguistics and Philosophy*, 7(2):135–154, 1984.
- C.L. Sidner. Plan parsing for intended response recognition in discourse. *Computational intelligence*, 1(1):1–10, 1985.

- K. VanLehn, C. Lynch, K. Schulze, J. A. Shapiro, R. H. Shelby, L. Taylor, D. J. Treacy, A. Weinstein, and M. C. Wintersgill. The Andes physics tutoring system: Lessons learned. *International Journal of Artificial Intelligence and Education*, 15(3), 2005.
- M. B. Vilain. Getting serious about parsing plans: A grammatical analysis of plan recognition. In *AAAI*, pages 190–197, 1990.
- R. Weida and D. Litman. Terminological reasoning with constraint networks and an application to plan recognition. In *Proc. of the 3rd Int. Conf. on the Principles of Knowledge Representation and Reasoning (KR-92)*, pages 282–293. Citeseer, 1992.
- R. Wilensky. Why John married Mary: Understanding stories involving recurring goals. *Cognitive Science*, 2(3):235–266, 1978.
- D. J. Yaron, M. Karabinos, A. Borek, B. McLaren, K. L. Evans, and G. Leinhardt. Tracking and supporting learning in open-ended activities involving a virtual lab simulation. In *The 238th American Chemical Society National Meeting, Washington, DC*, 2009.