

GSAT and Dynamic Backtracking

Matthew L. Ginsberg*
CIRL
1269 University of Oregon
Eugene, OR 97403
ginsberg@cirl.uoregon.edu

David A. McAllester†
AT&T Bell Laboratory
Murray Hill, New Jersey 07974
dmac@research.att.com

Abstract

There has been substantial recent interest in two new families of search techniques. One family consists of nonsystematic methods such as GSAT; the other contains systematic approaches that use a polynomial amount of justification information to prune the search space. This paper introduces a new technique that combines these two approaches. The algorithm allows substantial freedom of movement in the search space but enough information is retained to ensure the systematicity of the resulting analysis. The size of the justification database is guaranteed to be polynomial in the size of the problem in question.

1 INTRODUCTION

The past few years have seen rapid progress in the development of algorithms for solving constraint-satisfaction problems, or CSPs. CSPs arise naturally in subfields of AI from planning to vision, and examples include propositional theorem proving, map coloring and scheduling problems. The problems are difficult because they involve search; there is never a guarantee that (for example) a successful coloring of a portion of a large map can be extended to a coloring of the map in its entirety.

The algorithms developed recently have been of two types. *Systematic* algorithms determine whether a solution exists by searching the entire space. *Local* algorithms use hill-climbing techniques to find a solution quickly but are *nonsystematic* in that they search the entire space in only a probabilistic sense.

The empirical effectiveness of these nonsystematic algorithms appears to be a result of their ability to follow local gradients in the search space. Traditional

systematic procedures explore the space in a fixed order that is independent of local gradients; the fixed order makes following local gradients impossible but is needed to ensure that no node is examined twice and that the search remains systematic.

Dynamic backtracking [Ginsberg,1993] attempts to overcome this problem by retaining specific information about those portions of the search space that have been eliminated and then following local gradients in the remainder. Unlike previous algorithms that recorded such elimination information, such as dependency-directed backtracking [Stallman and Sussman,1977], dynamic backtracking is selective about the information it caches so that only a polynomial amount of memory is required. These earlier techniques cached a new result with every backtrack, using an amount of memory that was linear in the run time and thus exponential in the size of the problem being solved.

Unfortunately, neither dynamic nor dependency-directed backtracking (or any other known similar method) is truly effective at local maneuvering within the search space, since the basic underlying methodology remains simple chronological backtracking. New techniques are included to make the search more efficient, but an exponential number of nodes in the search space must still be examined before early choices can be retracted. No existing search technique is able to both move freely within the search space and keep track of what has been searched and what hasn't.

The second class of algorithms developed recently presume that freedom of movement is of greater importance than systematicity. Algorithms in this class achieve their freedom of movement by abandoning the conventional description of the search space as a tree of partial solutions, instead thinking of it as a space of total assignments of values to variables. Motion is permitted between any two assignments that differ on a single value, and a hill-climbing procedure is employed to try to minimize the number of constraints violated

*Supported by the Air Force Office of Scientific Research under contract 92-0693, by ARPA/Rome Labs under contracts numbers F30602-91-C-0036 and F30602-93-C-00031. This paper appeared in KR-94.

†Supported by ARPA under contract F33615-91-C-1788.

by the overall assignment. The best-known algorithms in this class are min-conflicts [Minton *et al.*,1990] and GSAT [Selman *et al.*,1992].

Min-conflicts has been applied to the scheduling domain specifically and used to schedule tasks on the Hubble space telescope. GSAT is restricted to Boolean satisfiability problems (where every variable is assigned simply true or false), and has led to remarkable progress in the solution of randomly generated problems of this type; its performance is reported [Selman and Kautz,1993, Selman *et al.*,1992, Selman *et al.*,1993] as surpassing that of other techniques such as simulated annealing [Kirkpatrick *et al.*,1982] and systematic techniques based on the Davis-Putnam procedure [Davis and Putnam,1960].

GSAT is not a panacea, however; there are many problems on which it performs fairly poorly. If a problem has no solution, for example, GSAT will never be able to report this with confidence. Even if a solution does exist, there appear to be at least two possible difficulties that GSAT may encounter.

First, the GSAT search space may contain so many local minima that it is not clear how GSAT can move so as to reduce the number of constraints violated by a given assignment. As an example, consider the CSP of generating crossword puzzles by filling words from a fixed dictionary into an empty frame [Ginsberg *et al.*,1990]. The constraints indicate that there must be no conflict in each of the squares; thus two words that begin on the same square must also begin with the same letter. In this domain, getting “close” is not necessarily any indication that the problem is nearly solved, since correcting a conflict at a single square may involve modifying much of the current solution. Konolige has recently reported that GSAT specifically has difficulty solving problems of this sort [Konolige,1994].

Second, GSAT does no forward propagation. In the crossword domain once again, selecting one word may well force the selection of a variety of subsequent words. In a Boolean satisfiability problem, assigning one variable the value true may cause an immediate cascade of values to be assigned to other variables via a technique known as *unit resolution*. It seems plausible that forward propagation will be more common on realistic problems than on randomly generated ones; the most difficult random problems appear to be tangles of closely related individual variables while naturally occurring problems tend to be tangles of sequences of related variables. Furthermore, it appears that GSAT’s performance degrades (relative to systematic approaches) as these sequences of variables arise [Crawford and Baker,1994].

Our aim in this paper is to describe a new search procedure that appears to combine the benefits of both of the earlier approaches; in some very loose sense, it can be thought of as a systematic version of GSAT.

The next three sections summarize the original dynamic backtracking algorithm [Ginsberg,1993], presenting it from the perspective of local search. The termination proof is omitted here but can be found in earlier papers [Ginsberg,1993, McAllester,1993]. Section 5 present a modification of dynamic backtracking called *partial-order dynamic backtracking*, or PDB. This algorithm builds on work of McAllester’s [McAllester,1993]. Partial-order dynamic backtracking provides greater flexibility in the allowed set of search directions while preserving systematicity and polynomial worst case space usage. Section 6 presents some empirical results comparing PDB with other well known algorithms on a class of “local” randomly generated 3-SAT problems. Concluding remarks are contained in Section 7.

2 CONSTRAINTS AND NOGOODS

We begin with a slightly nonstandard definition of a CSP.

Definition 2.1 *By a constraint satisfaction problem (I, V, κ) we will mean a finite set I of variables; for each $x \in I$, there is a finite set V_x of possible values for the variable x . κ is a set of constraints each of the form $\neg[(x_1 = v_1) \wedge \dots \wedge (x_k = v_k)]$ where each x_j is a variable in I and each v_j is an element of V_{x_j} . A solution to the CSP is an assignment P of values to variables that satisfies every constraint. For each variable x we require that $P(x) \in V_x$ and for each constraint $\neg[(x_1 = v_1) \wedge \dots \wedge (x_k = v_k)]$ we require that $P(x_i) \neq v_i$ for some x_i .*

By the size of a constraint-satisfaction problem (I, V, κ) , we will mean the product of the domain sizes of the various variables, $\prod_x |V_x|$.

The technical convenience of the above definition of a constraint will be clear shortly. For the moment, we merely note that the above description is clearly equivalent to the conventional one; rather than represent the constraints in terms of allowed value combinations for various variables, we write axioms that disallow specific value combinations one at a time. The size of a CSP is the number of possible assignments of values to variables.

Systematic algorithms attempting to find a solution to a CSP typically work with partial solutions that are then discovered to be inextensible or to violate the given constraints; when this happens, a backtrack occurs and the partial solution under consideration is modified. Such a procedure will, of course, need to

record information that guarantees that the same partial solution not be considered again as the search proceeds. This information might be recorded in the structure of the search itself; depth-first search with chronological backtracking is an example. More sophisticated methods maintain a database of some form indicating explicitly which choices have been eliminated and which have not. In this paper, we will use a database consisting of a set of *nogoods* [de Kleer,1986].

Definition 2.2 A nogood is an expression of the form

$$(x_1 = v_1) \wedge \dots \wedge (x_k = v_k) \rightarrow x \neq v \quad (1)$$

A nogood can be used to represent a constraint as an implication; (1) is logically equivalent to the constraint

$$\neg[(x_1 = v_1) \wedge \dots \wedge (x_k = v_k) \wedge (x = v)]$$

There are clearly many different ways of representing a given constraint as a nogood.

One special nogood is the *empty* nogood, which is tautologically false. We will denote the empty nogood by \perp ; if \perp can be derived from the given set of constraints, it follows that no solution exists for the problem being attempted.

The typical way in which new nogoods are obtained is by resolving together old ones. As an example, suppose we have derived the following:

$$\begin{aligned} (x = a) \wedge (y = b) &\rightarrow u \neq v_1 \\ (x = a) \wedge (z = c) &\rightarrow u \neq v_2 \\ (y = b) &\rightarrow u \neq v_3 \end{aligned}$$

where v_1, v_2 and v_3 are the only values in the domain of u . It follows that we can combine these nogoods to conclude that there is no solution with

$$(x = a) \wedge (y = b) \wedge (z = c) \quad (2)$$

Moving z to the conclusion of (2) gives us

$$(x = a) \wedge (y = b) \rightarrow z \neq c$$

In general, suppose we have a collection of nogoods of the form

$$x_{i1} = v_{i1} \wedge \dots \wedge x_{in_i} = v_{in_i} \rightarrow x \neq v_i$$

as i varies, where the same variable appears in the conclusions of all the nogoods. Suppose further that the antecedents all agree as to the value of the x_i 's, so that any time x_i appears in the antecedent of one of the nogoods, it is in a term $x_i = v_i$ for a fixed v_i . If the nogoods collectively eliminate all of the possible values

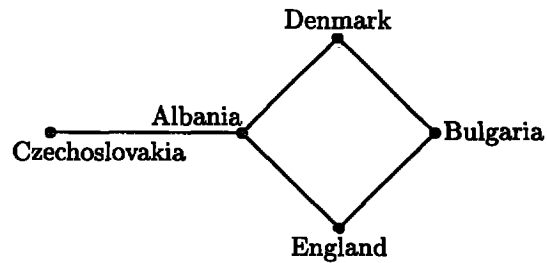


Figure 1: A small map-coloring problem

for x , we can conclude that $\bigwedge_j (x_j = v_j)$ is inconsistent; moving one specific x_k to the conclusion gives us

$$\bigwedge_{j \neq k} (x_j = v_j) \rightarrow x_k \neq v_k \quad (3)$$

As before, note the freedom in our choice of variable appearing in the conclusion of the nogood. Since the next step in our search algorithm will presumably satisfy (3) by changing the value for x_k , the selection of consequent variable corresponds to the choice of variable to “flip” in the terms used by GSAT or other hill-climbing algorithms.

As we have remarked, dynamic backtracking accumulates information in a set of nogoods. To see how this is done, consider the map coloring problem in Figure 1, repeated from [Ginsberg,1993]. The map consists of five countries: Albania, Bulgaria, Czechoslovakia, Denmark and England. We assume – wrongly – that the countries border each other as shown in the figure, where countries are denoted by nodes and border one another if and only if there is an arc connecting them.

In coloring the map, we can use the three colors red, green and blue. We will typically abbreviate the colors and country names to single letters in the obvious way. The following table gives a trace of how a conventional dependency-directed backtracking scheme might attack this problem; each row shows a state of the procedure in the middle of a backtrack step, after a new nogood has been identified but before colors are erased to reflect the new conclusion. The coloring that is about to be removed appears in boldface. The “drop”

column will be discussed shortly.

A	B	C	D	E	add	drop
r	g	r			$A = r \rightarrow C \neq r$	1
r	g	b	r		$A = r \rightarrow D \neq r$	2
r	g	b	g		$B = g \rightarrow D \neq g$	3
r	g	b	b	r	$A = r \rightarrow E \neq r$	4
r	g	b	b	g	$B = g \rightarrow E \neq g$	5
r	g	b	b	b	$D = b \rightarrow E \neq b$	6
r	g	b	b		$(A = r) \wedge (B = g) \rightarrow D \neq b$	7 6
r	g	b			$A = r \rightarrow B \neq g$	8 3,5,7

We begin by coloring Albania red and Bulgaria green, and then try to color Czechoslovakia red as well. Since this violates the constraint that Albania and Czechoslovakia be different colors, nogood (1) in the above table is produced.

We change Czechoslovakia's color to blue and then turn to Denmark. Since Denmark cannot be colored red or green, nogoods (2) and (3) appear; the only remaining color for Denmark is blue.

Unfortunately, having colored Denmark blue, we cannot color England. The three nogoods generated are (4), (5) and (6), and we can resolve these together because the three conclusions eliminate all of the possible colors for England. The result is that there is no solution with $(A = r) \wedge (B = g) \wedge (D = b)$, which we rewrite as (7) above. This can in turn be resolved with (2) and (3) to get (8), correctly indicating that the color of red for Albania is inconsistent with the choice of green for Bulgaria. The analysis can continue at this point to gradually determine that Bulgaria has to be red, Denmark can be green or blue, and England must then be the color not chosen for Denmark.

As we mentioned in the introduction, the problem with this approach is that the set Γ of nogoods grows monotonically, with a new nogood being added at every step. The number of nogoods stored therefore grows linearly with the run time and thus (presumably) exponentially with the size of the problem. A related problem is that it may become increasingly difficult to extend the partial solution P without violating one of the nogoods in Γ .

Dynamic backtracking deals with this by discarding nogoods when they become "irrelevant" in the sense that their antecedents no longer match the partial solution in question. In the example above, nogoods can be eliminated as indicated in the final column of the trace. When we derive (7), we remove (6) because Denmark is no longer colored blue. When we derive (8), we remove all of the nogoods with $B = g$ in their antecedents. Thus the only information we retain is that Albania's red color precludes red for Czechoslovakia,

Denmark and England (1, 2 and 4) and also green for Bulgaria (8).

3 DYNAMIC BACKTRACKING

Dynamic backtracking uses the set of nogoods to both record information about the portion of the search space that has been eliminated and to record the current partial assignment being considered by the procedure. The current partial assignment is encoded in the antecedents of the current nogood set. More formally:

Definition 3.1 An acceptable next assignment for a nogood set Γ is an assignment P satisfying every nogood in Γ and every antecedent of every such nogood. We will call a set of nogoods Γ acceptable if no two nogoods in Γ have the same conclusion and either $\perp \in \Gamma$ or there exists an acceptable next assignment for Γ .

If Γ is acceptable, the antecedents of the nogoods in Γ induce a partial assignment of values to variables; any acceptable next assignment must be an extension of this partial assignment. In the above table, for example, nogoods (1) through (6) encode the partial assignment given by $A = r$, $B = g$, and $D = b$. Nogoods (1) though (7) fail to encode a partial assignment because the seventh nogood is inconsistent with the partial assignment encoded in nogoods (1) through (6). This is why the sixth nogood is removed when the seventh nogood is added.

Procedure 3.2 (Dynamic backtracking) To solve a CSP:

$P :=$ any complete assignment of values to variables

$\Gamma := \emptyset$

until either P is a solution or $\perp \in \Gamma$:

$\gamma :=$ any constraint violated by P

$\Gamma := \text{simp}(\Gamma \cup \gamma)$

$P :=$ any acceptable next assignment for Γ

To simplify the discussion we assume a fixed total order on the variables. Versions of dynamic backtracking with dynamic rearrangement of the variable order can be found elsewhere [Ginsberg,1993, McAllester,1993]. Whenever a new nogood is added, the fixed variable ordering is used to select the variable that appears in the conclusion of the nogood – the latest variable always appears in the conclusion. The subroutine `simp` closes the set of nogoods under the resolution inference rule discussed in the previous section and removes all nogoods which have an antecedent $x = v$ such that $x \neq v$ appears in the conclusion of some other nogood. Without giving a detailed analysis, we note that simplification ensures that Γ remains acceptable. To prove termination we introduce the following notation:

Definition 3.3 For any acceptable Γ and variable x , we define the live domain of x to be those values v such that $x \neq v$ does not appear in the conclusion of any nogood in Γ . We will denote the size of the live domain of x by $|x|_\Gamma$, and will denote by $m(\Gamma)$ the tuple $\langle |x_1|_\Gamma, \dots, |x_n|_\Gamma \rangle$ where x_1, \dots, x_n are the variables in the CSP in their specified order.

Given an acceptable Γ , we define the size of Γ to be

$$\text{size}(\Gamma) = \prod_x |V_x| - \sum_x \left[(|V_x| - |x|_\Gamma) \prod_{x_i > x} |V_{x_i}| \right]$$

Informally, the size of Γ is the size of the remaining search space given the live domains for the variables and assuming that all information about x_i will be lost when we change the value for any variable $x_j < x_i$.

The following result is obvious:

Lemma 3.4 Suppose that Γ and Γ' are such that $m(\Gamma)$ is lexicographically less than $m(\Gamma')$. Then $\text{size}(\Gamma) < \text{size}(\Gamma')$. ■

The termination proof (which we do not repeat here) is based on the observation that every simplification lexicographically reduces $m(\Gamma)$. Assuming that $\Gamma = \emptyset$ initially, since

$$\text{size}(\emptyset) = \prod_x |V_x|$$

it follows that the running time of dynamic backtracking is bounded by the size of the problem being solved.

Proposition 3.5 Any acceptable set of nogoods can be stored in $o(n^2v)$ space where n is the number of variables and v is the maximum domain size of any single variable.

It is worth considering the behavior of Procedure 3.2 when applied to a CSP that is the union of two disjoint CSPs that do not share variables or constraints. If each of the two subproblems is unsatisfiable and the variable ordering interleaves the variables of the two subproblems, a classical backtracking search will take time proportional to the product of the times required to search each assignment space separately.¹ In contrast, Procedure 3.2 works on the two problems independently, and the time taken to solve the union of problems is therefore the sum of the times needed for the individual subproblems. It follows that Procedure 3.2 is fundamentally different from classical backtracking or backjumping procedures; Procedure 3.2 is in fact what has been called a *polynomial space aggressive backtracking procedure* [McAllester,1993].

¹This observation remains true even if backjumping techniques are used.

4 DYNAMIC BACKTRACKING AS LOCAL SEARCH

Before proceeding, let us highlight the obvious similarities between Procedure 3.2 and Selman's description of GSAT [Selman *et al.*,1992]:

Procedure 4.1 (GSAT) To solve a CSP:

```

for i := 1 to MAX-TRIES
  P := a randomly generated truth assignment
  for j := 1 to MAX-FLIPS
    if P is a solution, then return it
    else flip any variable in P that results in
         the greatest decrease in the number
         of unsatisfied clauses
  end if
end for
return failure
    
```

The inner loop of the above procedure makes a local move in the search space in a direction consistent with the goal of satisfying a maximum number of clauses; we will say that GSAT follows the local gradient of a "maxsat" objective function. But local search can get stuck in local minima; the outer loop provides a partial escape by giving the procedure several independent chances to find a solution.

Like GSAT, dynamic backtracking examines a sequence of total assignments. Initially, dynamic backtracking has considerable freedom in selecting the next assignment; in many cases, it can update the total assignment in a manner identical to GSAT. The nogood set ultimately both constrains the allowed directions of motion and forces the procedure to search systematically. Dynamic backtracking cannot get stuck in local minima.

Both systematicity and the ability to follow local gradients are desirable. The observations of the previous paragraphs, however, indicate that these two properties are in conflict – systematic enumeration of the search space appears incompatible with gradient descent. To better understand the interaction of systematicity and local gradients, we need to examine more closely the structure of the nogoods used in dynamic backtracking.

We have already discussed the fact that a single constraint can be represented as a nogood in a variety of ways. For example, the constraint $\neg(A = r \wedge B = g)$ can be represented either as $A = r \rightarrow B \neq g$ or as $B = g \rightarrow A \neq r$. Although these nogoods capture the same information, they behave differently in the dynamic backtracking procedure because they encode different partial truth assignments and represent dif-

ferent choices of variable ordering. In particular, the set of acceptable next assignments for $A = r \rightarrow B \neq g$ is quite different from the set of acceptable next assignments for $B = g \rightarrow A \neq r$. In the former case an acceptable assignment must satisfy $A = r$; in the latter case, $B = g$ must hold. Intuitively, the former nogood corresponds to changing the value of B while the latter nogood corresponds to changing that of A . The manner in which we represent the constraint $\neg(A = r \wedge B = g)$ influences the direction in which the search is allowed to proceed. In Procedure 3.2, the choice of representation is forced by the need to respect the fixed variable ordering and to change the latest variable in the constraint.² Similar restrictions exist in the original presentation of dynamic backtracking itself [Ginsberg,1993].

5 PARTIAL-ORDER DYNAMIC BACKTRACKING

Partial-order dynamic backtracking [McAllester,1993] replaces the fixed variable order with a *partial* order that is dynamically modified during the search. When a new nogood is added, this partial ordering need not fix a unique representation – there can be considerable choice in the selection of the variable to appear in the conclusion of the nogood. This leads to freedom in the selection of the variable whose value is to be changed, thereby allowing greater flexibility in the directions that the procedure can take while traversing the search space. The locally optimal gradient followed by GSAT can be adhered to more often. The partial order on variables is represented by a set of ordering constraints called *safety conditions*.

Definition 5.1 *A safety condition is an assertion of the form $x < y$ where x and y are variables. Given a set S of safety conditions, we will denote by \leq_S the transitive closure of $<$, and will require that \leq_S be antisymmetric. We will write $x <_S y$ to mean that $x \leq_S y$ and $y \not\leq_S x$.*

In other words, $x \leq y$ if there is some (possibly empty) sequence of safety conditions

$$x < z_1 < \dots < z_n < y$$

The requirement of antisymmetry means simply that there are no two distinct x and y for which $x \leq y$ and $y \leq x$; in other words, \leq_S has no “loops” and is a partial order on the variables.

Definition 5.2 *For a nogood γ , we will denote by S_γ the set of all safety conditions $x < y$ such that x is in*

²Note, however, that there is still considerable freedom in the choice of the constraint itself. A total assignment usually violates many different constraints.

the antecedent of γ and y is the variable in its conclusion.

Informally, we require variables in the antecedent of nogoods to precede the variables in their conclusions, since the antecedent variables have been used to constrain the live domains of the conclusions.

The state of the partial order dynamic backtracking procedure is represented by a pair $\langle \Gamma, S \rangle$ consisting of a set of nogoods and a set of safety conditions. In many cases, we will be interested in only the ordering information about variables that can precede a fixed variable x . To discard the rest of the ordering information, we discard all of the safety conditions involving any variable y that follows x , and then record only that y does indeed follow x . Somewhat more formally:

Definition 5.3 *For any set S of safety conditions and variable x , we define the weakening of S at x , to be denoted $W(S, x)$, to be the set of safety conditions given by removing from S all safety conditions of the form $z < y$ where $x <_S y$ and then adding the safety condition $x < y$ for all such y .*

The set $W(S, x)$ is a weakening of S in the sense that every total ordering consistent with S is also consistent with $W(S, x)$. However $W(S, x)$ usually admits more total orderings than S does; for example, if S specifies a total order then $W(S, x)$ allows any order which agrees with S up to and including the variable x . In general, we have the following:

Lemma 5.4 *For any set S of safety conditions, variable x , and total order $<$ consistent with the safety conditions in $W(S, x)$, there exists a total order consistent with S that agrees with $<$ through x .*

We now state the PDB procedure.

Procedure 5.5 To solve a CSP:

$P :=$ any complete assignment of values to variables

$\Gamma := \emptyset$

$S := \emptyset$

until either P is a solution or $\perp \in \Gamma$:

$\gamma :=$ a constraint violated by P

$\langle \Gamma, S \rangle := \text{simp}(\Gamma, S, \gamma)$

$P :=$ any acceptable next assignment for Γ

Procedure 5.6 To compute $\text{simp}(\Gamma, S, \gamma)$:
 select the conclusion x of γ so that $S \cup S_\gamma$ is acyclic
 $\Gamma := \Gamma \cup \{\gamma\}$
 $S := W(S \cup S_\gamma, x)$
 remove from Γ each nogood with x in its antecedent
if the conclusions of nogoods in Γ rule out all possible values for x **then**
 $\rho :=$ the result of resolving all nogoods in Γ with x in their conclusion
 $\langle \Gamma, S \rangle := \text{simp}(\Gamma, S, \rho)$
end if
return $\langle \Gamma, S \rangle$

The above simplification procedure maintains the invariant that Γ be acceptable and S be acyclic; in addition, the time needed for a single call to simp appears to grow significantly sublinearly with the size of the problem in question (see Section 6).

Theorem 5.7 *Procedure 5.5 terminates. The number of calls to simp is bounded by the size of the problem being solved.*

As an example, suppose that we return to our map-coloring problem. We begin by coloring all of the countries red except Bulgaria, which is green. The following table shows the total assignment that existed at the moment each new nogood was generated.

A	B	C	D	E	add	drop
r	g	r	r	r	$C = r \rightarrow A \neq r$	1
b	g	r	r	r	$D = r \rightarrow E \neq r$	2
b	g	r	r	g	$B = g \rightarrow E \neq g$	3
b	g	r	r	b	$A = b \rightarrow E \neq b$	4
					$(A = b) \wedge (B = g) \rightarrow D \neq r$	5 2
					$D < E$	6
b	g	r	g	r	$B = g \rightarrow D \neq g$	7
b	g	r	b	r	$A = b \rightarrow D \neq b$	8
					$A = b \rightarrow B \neq g$	9 3,5,7
					$B < E$	10 6
					$B < D$	11

The initial coloring violates a variety of constraints; suppose that we choose to work on one with Albania in its conclusion because Albania is involved in three violated constraints. We choose $C = r \rightarrow A \neq r$ specifically, and add it as (1) above.

We now modify Albania to be blue. The only constraint violated is that Denmark and England be different colors, so we add (2) to Γ . This suggests that we change the color for England; we try green, but this conflicts with Bulgaria. If we write the new nogood as $E = g \rightarrow B \neq g$, we will change Bulgaria to blue and be done. In the table above, however, we have made

the less optimal choice (3), changing the coloring for England again.

We are now forced to color England blue. This conflicts with Albania, and we continue to leave England in the conclusion of the nogood as we add (4). This nogood resolves with (2) and (3) to produce (5), where we have once again made the worst choice and put D in the conclusion. We add this nogood to Γ and remove nogood (2), which is the only nogood with D in its antecedent. In (6) we add a safety condition indicating that D must continue to precede E . (This safety condition has been present since nogood (2) was discovered, but we have not indicated it explicitly until the original nogood was dropped from the database.)

We next change Denmark to green; England is forced to be red once again. But now Bulgaria and Denmark are both green; we have to write this new nogood (7) with Denmark in the conclusion because of the ordering implied by nogood (5) above. Changing Denmark to blue conflicts with Albania (8), which we have to write as $A = b \rightarrow D \neq b$. This new nogood resolves with (5) and (7) to produce (9).

We drop (3), (5) and (7) because they involve $B = g$, and introduce the two safety conditions (10) and (11). Since E follows B , we drop the safety condition $E < D$. At this point, we are finally forced to change the color for Bulgaria and the search continues.

It is important to note that the added flexibility of PDB over dynamic backtracking arises from the flexibility in the first step of the simplification procedure where the conclusion of the new nogood is selected. This selection corresponds to a selection of a variable whose value is to be changed.

As with the procedure in the previous section, when given a CSP that is a union of disjoint CSPs the above procedure will treat the two subproblems independently. The total running time remains the sum of the times required for the subproblems.

6 EXPERIMENTAL RESULTS

In this section, we present preliminary results regarding the implemented effectiveness of the procedure we have described. We compared a search engine based on this procedure with two others, TABLEAU [Crawford and Auton,1995] and WSAT, or "walk-sat" [Selman *et al.*,1993]. TABLEAU is an efficient implementation of the Davis-Putnam algorithm and is systematic; WSAT is a modification to GSAT and is not. We used WSAT instead of GSAT because WSAT is more effective on a fairly wide range of problem distributions [Selman *et al.*,1993].

The experimental data was not collected using the random 3-SAT problems that have been the

target of much recent investigation, since there is growing evidence that these problems are not representative of the difficulties encountered in practice [Crawford and Baker,1994]. Instead, we generated our problems so that the clauses they contain involve groups of locally connected variables as opposed to variables selected at random.

Somewhat more specifically, we filled an $n \times n$ square grid with variables, and then required that the three variables appearing in any single clause be neighbors in this grid. We believe that the qualitative properties of the results reported here hold for a wide class of distributions where variables are given spatial locations and clauses are required to be local.

The experiments were performed at the crossover point where approximately half of the instances generated could be expected to be satisfiable, since this appears to be where the most difficult problems lie [Crawford and Auton,1995]. Note that not all instances at the crossover point are hard; as an example, the local variable interactions in these problems can lead to short resolution proofs that no solution exists in unsatisfiable cases. This is in sharp contrast with random 3-SAT problems (where no short proofs appear to exist in general, and it can even be shown that proof lengths are growing exponentially on average [Chvátal and Szemerédi,1988]). Realistic problems may often have short proof paths: A particular scheduling problem may be unsatisfiable simply because there is no way to schedule a specific resource as opposed to because of global issues involving the problem in its entirety. Satisfiability problems arising in VLSI circuit design can also be expected to have locality properties similar to those we have described.

The problems involved 25, 100, 225, 400 and 625 variables. For each size, we generated 100 satisfiable and 100 unsatisfiable instances and then executed the three procedures to measure their performance. (WSAT was not tested on the unsatisfiable instances.) For WSAT, we measured the number of times specific variable values were flipped. For PDB, we measured the number of top-level calls to Procedure 5.6. For TABLEAU, we measured the number of choice nodes expanded. WSAT and PDB were limited to 100,000 flips; TABLEAU was limited to a running time of 150 seconds.

The results for the satisfiable problems were as follows. For TABLEAU, we give the node count for successful runs only; we also indicate parenthetically what fraction of the problems were solved given the computational resource limitations. (WSAT and PDB successfully solved all instances.)

Variables	PDB	WSAT	TABLEAU
25	35	89	9 (1.0)
100	210	877	255 (.98)
225	434	1626	504 (.70)
400	731	2737	856 (.70)
625	816	3121	502 (.68)

For the unsatisfiable instances, the results were:

Variables	PDB	TABLEAU
25	122	8 (1.0)
100	509	1779 (1.0)
225	988	5682 (.38)
400	1090	558 (.11)
625	1204	114 (.06)

The times required for PDB and WSAT appear to be growing comparably, although only PDB is able to solve the unsatisfiable instances. The eventual *decrease* in the average time needed by TABLEAU is because it is only managing to solve the easiest instances in each class. This causes TABLEAU to become almost completely ineffective in the unsatisfiable case and only partially effective in the satisfiable case. Even where it does succeed on large problems, TABLEAU's run time is greater than that of the other two methods.

Finally, we collected data on the time needed for each top-level call to `simp` in partial-order dynamic backtracking. As a function of the number of variables in the problem, this was:

Number of variables	PDB (msec)	WSAT (msec)
25	3.9	0.5
100	5.3	0.3
225	6.7	0.6
400	7.0	0.7
625	8.4	1.4

All times were measured on a Sparc 10/40 running unoptimized Allegro Common Lisp. An efficient C implementation could expect to improve either method by approximately an order of magnitude. As mentioned in Section 5, the time per flip is growing sublinearly with the number of variables in question.

7 CONCLUSION AND FUTURE WORK

Our aim in this paper has been to make a primarily theoretical contribution, describing a new class of constraint-satisfaction algorithms that appear to combine many of the advantages of previous systematic and nonsystematic approaches. Since our focus has been on a description of the algorithms, there is obviously much that remains to be done.

First, of course, the procedures must be tested on a variety of problems, both synthetic and nat-

urally occurring; the results reported in Section 6 only scratch the surface. It is especially important that realistic problems be included in any experimental evaluation of these ideas, since these problems are likely to have performance profiles substantially different from those of randomly generated problems [Crawford and Baker,1994]. The experiments of the previous section need to be extended to include unit resolution.

Finally, we have left completely untouched the question of how the flexibility of Procedure 5.6 is to be exploited. Given a group of violated constraints, which should we pick to add to Γ ? Which variable should be in the conclusion of the constraint? These choices correspond to choice of backtrack strategy in a more conventional setting, and it will be important to understand them in this setting as well.

Acknowledgement

We would like to thank Jimi Crawford, Ari Jónsson, Bart Selman and the members of CIRL for taking the time to discuss these ideas with us. Crawford especially contributed to the development of Procedure 5.5.

References

- [Chvátal and Szemerédi,1988] V. Chvátal and E. Szemerédi. Many hard examples for resolution. *JACM*, 35:759–768, 1988.
- [Crawford and Auton,1995] James M. Crawford and Larry D. Auton. Experimental results on the crossover point in random 3sat. *Artificial Intelligence*, 1995. To appear.
- [Crawford and Baker,1994] James M. Crawford and Andrew B. Baker. Experimental results on the application of satisfiability algorithms to scheduling problems. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, 1994.
- [Davis and Putnam,1960] M. Davis and H. Putnam. A computing procedure for quantification theory. *J. Assoc. Comput. Mach.*, 7:201–215, 1960.
- [de Kleer,1986] Johan de Kleer. An assumption-based truth maintenance system. *Artificial Intelligence*, 28:127–162, 1986.
- [Ginsberg et al.,1990] Matthew L. Ginsberg, Michael Frank, Michael P. Halpin, and Mark C. Torrance. Search lessons learned from crossword puzzles. In *Proceedings of the Eighth National Conference on Artificial Intelligence*, pages 210–215, 1990.
- [Ginsberg,1993] Matthew L. Ginsberg. Dynamic backtracking. *Journal of Artificial Intelligence Research*, 1:25–46, 1993.
- [Kirkpatrick et al.,1982] S. Kirkpatrick, C.D. Gelatt, and M.P. Vecchi. Optimization by simulated annealing. *Science*, 220:671–680, 1982.
- [Konolige,1994] Kurt Konolige. Easy to be hard: Difficult problems for greedy algorithms. In *Proceedings of the Fourth International Conference on Principles of Knowledge Representation and Reasoning*, Bonn, Germany, 1994.
- [McAllester,1993] David A. McAllester. Partial order backtracking. <ftp.ai.mit.edu:/pub/dam/dynamic.ps>, 1993.
- [Minton et al.,1990] Steven Minton, Mark D. Johnston, Andrew B. Philips, and Philip Laird. Solving large-scale constraint satisfaction and scheduling problems using a heuristic repair method. In *Proceedings of the Eighth National Conference on Artificial Intelligence*, pages 17–24, 1990.
- [Selman and Kautz,1993] Bart Selman and Henry Kautz. Domain-independent extensions to GSAT: Solving large structured satisfiability problems. In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence*, pages 290–295, 1993.
- [Selman et al.,1992] Bart Selman, Hector Levesque, and David Mitchell. A new method for solving hard satisfiability problems. In *Proceedings of the Tenth National Conference on Artificial Intelligence*, pages 440–446, 1992.
- [Selman et al.,1993] Bart Selman, Henry A. Kautz, and Bram Cohen. Local search strategies for satisfiability testing. In *Proceedings 1993 DIMACS Workshop on Maximum Clique, Graph Coloring, and Satisfiability*, 1993.
- [Stallman and Sussman,1977] R. M. Stallman and G. J. Sussman. Forward reasoning and dependency-directed backtracking in a system for computer-aided circuit analysis. *Artificial Intelligence*, 9(2):135–196, 1977.