



Centrum voor Wiskunde en Informatica

REPORT*RAPPORT*

A Two Phase Algorithm for Solving a Class of Hard Satisfiability Problems

J.P. Warners, H. van Maaren

Software Engineering (SEN)

SEN-R9802 April 30, 1998

Report SEN-R9802
ISSN 1386-369X

CWI
P.O. Box 94079
1090 GB Amsterdam
The Netherlands

CWI is the National Research Institute for Mathematics and Computer Science. CWI is part of the Stichting Mathematisch Centrum (SMC), the Dutch foundation for promotion of mathematics and computer science and their applications.

SMC is sponsored by the Netherlands Organization for Scientific Research (NWO). CWI is a member of ERCIM, the European Research Consortium for Informatics and Mathematics.

Copyright © Stichting Mathematisch Centrum
P.O. Box 94079, 1090 GB Amsterdam (NL)
Kruislaan 413, 1098 SJ Amsterdam (NL)
Telephone +31 20 592 9333
Telefax +31 20 592 4199

A Two Phase Algorithm for Solving a Class of Hard Satisfiability Problems

Joost P. Warners^{a,b} Hans van Maaren^b

^a CWI

P.O. Box 94079, 1090 GB, Amsterdam, The Netherlands

e-mail: Joost.Warners@cwi.nl

^b Department of Technical Mathematics and Informatics, Faculty of Information Technology and Systems
Delft University of Technology, P.O. Box 5031, 2600 GA Delft, The Netherlands

ABSTRACT

The DIMACS suite of satisfiability (SAT) benchmarks contains a set of instances that are very hard for existing algorithms. These instances arise from learning the parity function on 32 bits. In this paper we develop a two phase algorithm that is capable of solving these instances. In the first phase, a polynomially solvable subproblem is identified and solved. Using the solution to this problem, we can considerably restrict the size of the search-space in the second phase of the algorithm, which is an extension of the well-known Davis–Putnam–Loveland algorithm for SAT problems. We conclude with reporting on our computational results on the parity instances.

1991 Mathematics Subject Classification: 03B05, 68Q20, 68T20, 90C27

1991 Computing Reviews Classification System: F.2, G.2

Keywords and Phrases: Satisfiability, polynomial representation, Davis–Putnam algorithm

Note: The first author is supported by the Dutch Organization for Scientific Research (NWO) under grant SION 612-33-001. Research carried out under SEN 2.3

1. INTRODUCTION

In a recent paper by Selman et al. [8] ten challenges in propositional reasoning are formulated. One of these is to develop an efficient algorithm for solving instances arising from the parity learning problem on 32 bits [2]. Several instances of this problem are available in the DIMACS suite of SAT benchmarks [5]. None of the currently known algorithms appear to be capable of solving these instances in reasonable time. Incomplete algorithms do not succeed in finding models, while it seems that for systematic search procedures the search-space is too large [8].

We develop a two-phase algorithm for the parity problems that is capable of finding models in less than ten minutes. In the first phase of the algorithm a polynomially solvable subproblem is isolated and solved. The subproblem has a *balanced polynomial representation* [9], and can be shown to be equivalent to a formula that is a conjunction of (nested) equivalencies (CoE). Its solution allows us to reduce the search-space in the second phase considerably. In that phase we apply a DPL-type algorithm to a conjunction of a formula in Conjunctive Normal Form (CNF) and a CoE formula.

This paper is organized as follows. In the next section we discuss the necessary preliminaries. Subsequently, we introduce the concept of balanced polynomial representations (BPR) and show

that a formula with BPR is equivalent to a CoE formula. We briefly review a polynomial-time algorithm for CoE formulas. Section 4 is concerned with the recognition of CoE subformulas, and in Section 5 we extend the DPL algorithm to solve conjunctions of CNF and CoE formulas. We conclude with computational results.

2. PRELIMINARIES AND NOTATION

A propositional formula Φ in conjunctive normal form (CNF) is the conjunction of n clauses, where each clause is a disjunction of literals $(\neg)p_i$. Each literal is an *atomic proposition* (or *variable*) or its negation (\neg) . Let m be the number of atomic propositions. Thus each clause \mathbf{C}_k is of the form

$$\mathbf{C}_k = \bigvee_{i \in I_k} p_i \vee \bigvee_{j \in J_k} \neg p_j,$$

with $I_k, J_k \subseteq \{1, \dots, m\}$ disjoint. The satisfiability problem of propositional logic is to assign truth values to the variables, such that each clause evaluates to true (i.e. one of its literals is true) and so the whole formula evaluates to true, or it must be proved that no such assignment exists.

We define the matrix $A \in \mathbb{R}^{n \times m}$ to be the *clause-variable* matrix. Each row corresponds to a clause and each column is associated with a variable. It holds that $a_{ki} = 1$ if $i \in I_k$, $a_{ki} = -1$ if $i \in J_k$, while $a_{ki} = 0$ for any $i \notin I_k \cup J_k$. Note that, associating a $\{-1, 1\}$ variable x_i with each proposition letter p_i , the integer linear programming formulation of the satisfiability problem can be stated as finding a vector $x \in \{-1, 1\}^m$ such that $Ax \geq b$, where $b \in \mathbb{R}^n$, with $b_k = 2 - |I_k \cup J_k|$.

Now let us derive a different formulation of SAT problems, based on a multiplicative rather than additive representation of clauses. Formulations of this type have been used by Gu [4] to obtain effective approximation algorithms for large-scale satisfiability problems.

A clause \mathbf{C}_k is satisfied, if and only if $x \in \{-1, 1\}^m$ satisfies

$$P_k(x) = \prod_{i \in I_k} (1 - x_i) \prod_{j \in J_k} (1 + x_j) = \prod_{i=1}^m (1 - a_{ki}x_i) = 0. \quad (2.1)$$

Observe that $P_k(x)$ remains a valid representation of clause \mathbf{C}_k when multiplying it with a (strictly) positive weight w_k . Let $M = \{1, \dots, m\}$. In general, $x \in \{-1, 1\}^m$ is a satisfiable assignment of a formula Φ , if and only if

$$\mathcal{P}(x) = \sum_{k=1}^n w_k P_k(x) = \sum_{k=1}^n w_k + \sum_{I \subseteq M} (-1)^{|I|} \sum_{k=1}^n w_k \prod_{i \in I} a_{ki} x_i = 0,$$

where in principal I runs through all possible subsets of M ($I \neq \emptyset$) and w is a strictly positive weight vector. Note that the number of subsets that has to be taken into account can be restricted substantially, since in fact only subsets $I \subseteq M$ for which $I \subseteq I_k \cup J_k$ for some $k = 1, \dots, n$ need to be considered. In general, for a clause with length ℓ , $2^\ell - 1$ coefficients

need to be computed.

We use the notation

$$c_I = (-1)^{|I|} \sum_{k=1}^n w_k \prod_{i \in I} a_{ki}, \quad (2.2)$$

where $I \subseteq M = \{1, \dots, m\}$. The satisfiability problem has the following *polynomial representation*.

$$(PR) \quad \text{find } x \in \{-1, 1\}^m \text{ such that } \mathcal{P}(x) = \sum_{k=1}^n w_k + \sum_{I \subseteq M} c_I \prod_{i \in I} x_i = 0.$$

Observe that by construction $\mathcal{P}(x) \geq 0$ for any $x \in \{-1, 1\}^m$. Strict inequality implies that the corresponding CNF formula is unsatisfiable.

In this paper we also make use of propositional formulas in conjunction of equivalencies form (CoEs). Such formulas are solvable in polynomial time [9], as opposed to formulas in CNF which are in general NP-complete [1]. In the next section we briefly review a polynomial-time algorithm for CoE formulas.

A CoE formula is a conjunction of *equivalency-clauses*. An equivalency-clause \mathbf{Q}_k is defined as a (nested) equivalency of literals or its negation. We denote this as

$$\mathbf{Q}_k = [\neg]_k \bigcup_{i \in I_k} p_i. \quad (2.3)$$

Observe that the polynomial representation of \mathbf{Q}_k is very short:

$$Q_k(x) = \delta_k \prod_{i \in I_k} x_i = 1, \quad (2.4)$$

where $\delta_k = -1$ if the negation operator is present in (2.3), otherwise $\delta_k = 1$. This representation is obtained by directly considering (2.3); an equivalent representation is obtained by first translating \mathbf{Q}_k to CNF, and then summing the $2^{|I_k|-1}$ associated polynomial representations (2.1). Conversely, it is easy to see that to any equation of type (2.4) an equivalency-clause is associated. For example, if $\mathbf{Q}_k = \neg(p_1 \leftrightarrow p_4 \leftrightarrow p_8)$, then $Q_k(x) = -x_1 x_4 x_8 = 1$ and vice versa.

3. BALANCED POLYNOMIAL REPRESENTATIONS

In this section we discuss a notion of balancedness for SAT formulas, based on the polynomial representation (PR). The notions discussed here were earlier introduced in [9]. Let us start with a definition.

Definition 1 Consider the polynomial representation (PR). We call the polynomial function $\mathcal{P}(x)$ balanced if

$$\sum_{I \subseteq M} |c_I| = \sum_{k=1}^n w_k.$$

Furthermore, $\mathcal{P}(x)$ is called (strictly) positive if

$$\sum_{I \subseteq M} |c_I| < \sum_{k=1}^n w_k.$$

Now assume we are given a SAT formula Φ and its polynomial representation (PR). If $\mathcal{P}(x)$ is balanced, we say that Φ has a *balanced polynomial representation (BPR)*. Similarly, if $\mathcal{P}(x)$ is positive, we say that Φ has a *positive polynomial representation (PPR)*. In the latter case Φ is unsatisfiable [9].

We have the following lemma.

Lemma 1 *If Φ has a balanced polynomial representation, it is equivalent to a CoE formula.*

Proof: Observe that if $\mathcal{P}(x)$ is balanced, then for any feasible vector $x \in \{-1, 1\}^m$ it must hold that

$$c_I \prod_{i \in I} x_i = -|c_I|,$$

for all $I \subseteq M$. This implies that we may set c_I to $\text{sgn}(c_I)$, thus obtaining an equation of the form (2.4). \square

Let us now review a polynomial time algorithm for solving CoE formulas, which (implicitly) yields all satisfiable solutions. We only give the outline here, for a more detailed description the reader is referred to [9].

Consider an equivalency-clause \mathbf{Q}_k and its polynomial representation $Q_k(x)$ (2.4). Obviously, for any feasible solution $x \in \{-1, 1\}^m$ it holds that

$$x_j = \delta_k \prod_{i \in I_k \setminus j} x_i, \quad \text{for all } j \in I_k.$$

Choosing an index $j \in I_k$ we can substitute the above expression in all equivalency-clauses \mathbf{Q}_l ($l \neq k$) in which x_j occurs, using that $x_i^2 = 1$. Thus all but one occurrence of x_j are eliminated. Now the algorithm runs as follows. We initialize the set $\mathcal{I} = \{x_1, \dots, x_m\}$, the set of *independent* variables. We loop through the equivalency-clauses once, choosing a variable x_j in each one to eliminate from all other equivalency-clauses. Subsequently we remove x_j from \mathcal{I} , and call it a *dependent* variable. Thus we end up with a set of equivalency clauses, for which all satisfiable assignments can be constructed by assigning all possible combinations of truth values to the independent variables. The values of the dependent variables are uniquely determined by an assignment to the independent variables. Note that during the elimination process the equality $\delta_k = 1$ might be derived for a δ_k that is equal to -1 ; this implies that the formula is contradictory. Here is a small example.

Example. A balanced polynomial representation is given by

$$\mathcal{P}(x) = 7 - 2x_1x_2x_3 + x_1x_3x_5 - 3x_2x_4x_5 - x_1x_4x_5.$$

An equivalent representation is

$$\mathcal{P}^*(x) = 4 + x_1x_5 + x_2x_5 + x_4 - x_3,$$

with $\mathcal{I} = \{x_5\}$. Thus two distinct solutions can be constructed. \square

If a formula has a CoE *subformula*, solving this first may be of help in solving the full formula, since it allows us to take dependencies into account in a systematic way. When solving the full formula the search can possibly be restricted to the independent variables. Moreover, the CoE subformula might be a contradiction, implying that the full formula is also unsatisfiable.

4. POLYNOMIAL TIME RECOGNITION OF COE SUBFORMULAS

Let us now address the problem of recognizing a CoE subformula. We can make use of a linear programming (LP) formulation to find a CoE subformula of maximal weight. Since the construction of the LP can be done in polynomial time (assuming that the maximal clause length is bounded and fixed), and LP problems are polynomially solvable [6], the recognition problem can be solved in polynomial time.

In the formulation the weights w_k occurring in the polynomial representation (PR) are the main decision variables. Essentially, we want to find a set of nonnegative weights w_k and a *slack* $s \geq 0$ such that (see Definition 1 and equation (2.2)),

$$\sum_{I \subseteq M} \left| \sum_{k=1}^n \left(\prod_{i \in I} a_{ki} \right) w_k \right| + s = \sum_{k=1}^n w_k. \quad (4.1)$$

We allow the weights to be equal to zero; if $w_k = 0$ for some k , this implies that clause k is not in the subformula, while if $w_k > 0$ clause k is in the subformula. Our first goal should be to find a solution with s strictly positive (since then the associated subformula has PPR and is unsatisfiable); if no such solution exists, the goal is to identify a subformula of maximal weight with BPR. To check whether solutions with the desired properties exist, we first solve an LP with the objective of maximizing s , and if the optimal value of this LP is equal to zero, a second LP must be solved with the objective to maximize the sum of the weights. Consider the following LP.

$$\begin{aligned}
 (LP) \quad & \max \quad \alpha s + \beta \sum_{k=1}^n w_k \\
 & \text{s.t.} \quad \sum_{I \subseteq M} (z_I^+ + z_I^-) - \sum_{k=1}^n w_k + s = 0, \\
 & \quad \sum_{k=1}^n \left(\prod_{i \in I} a_{ki} \right) w_k - z_I^+ + z_I^- = 0, \quad I \subseteq M, \\
 & \quad 0 \leq w_k \leq 1, \quad 1 \leq k \leq n, \\
 & \quad z_I^+, z_I^- \geq 0, \quad I \subseteq M, \\
 & \quad s \geq 0.
 \end{aligned}$$

The two separate LPs are obtained by setting $\beta = 0$ and $\alpha = 0$, $s = 0$ respectively. The first constraint evaluates expression (4.1) and in the subsequent set of constraints the c_I are computed (see (2.2)). The auxiliary variables z_I^+ and z_I^- associated with the (nonempty) set I are used to eliminate the absolute values in (4.1) in the usual way. For a formula in which the clauses have a maximum length ℓ , the numbers of variables and constraints are bounded by $(2^{\ell+1} - 1)n + 1$ and $(2^\ell - 1)n + 1$ respectively.

Note that a subformula of maximal weight is not guaranteed to be a subformula of maximal *size*, although in most cases these will coincide. In this respect using an interior point method for solving (LP) is better than the simplex method, since an IPM yields an optimal solution with a maximal number of nonzero variables.

In practice, heuristics that look for particular structures may often succeed in identifying CoE subformulas. Indeed, for the parity formulas solved in this paper such heuristics suffice [9]. However, if a subformula is ‘well hidden’, or does not conform certain standard structures, using the LP approach described above will succeed in identifying it, whereas the heuristic methods are likely to fail.

Observe that if the optimal value of the first LP is equal to zero, no subformula with PPR exists. Obviously, the existence of a subformula with PPR is merely a sufficient condition for a formula to be contradictory. If the optimal value of the second LP equals zero, no CoE subformula exists. For random instances this will usually be the case. On the other hand, instances that stem from some practical application often have a lot of structure that can be utilized via this LP approach.

5. A DPL ALGORITHM FOR SOLVING MIXED CNF/CoE FORMULAS

One of the best known exact algorithms for solving CNF formulas is the variant of the Davis–Putnam algorithm [3] introduced by Loveland [7], which is known as the Davis–Putnam–Loveland (DPL) algorithm. The DPL–algorithm implicitly enumerates all 2^m distinct solutions. by setting up a binary search tree. We can easily extend this algorithm to solve conjunctions of CNF and CoE formulas. In figure 1 the extension of the algorithm is summarized.

Let us look a bit more closely at the algorithm. First we consider the unit resolution phase. When a unit literal is propagated through the formula, some clauses become true, while others reduce in length by one. For equivalency–clauses it holds that each in which the current unit literal occurs simply reduces in length by one. As usual, unit resolution is applied until no unit clauses remain, where it is noted that an equivalency clause of length one can be regarded as a unit clause in the usual sense. After the unit resolution phase it is checked whether the current formula can be declared either satisfiable or contradictory. If not, a *branching* or *splitting* variable l is chosen in some pre–specified way and the DPL procedure is recursively called with this variable set to true and false respectively.

```

procedure DPL ( $\Phi = \Phi_{CNF} \cup \Phi_{CoE}$ , depth);
   $\Phi := \text{unit\_resolution}(\Phi)$ ;
  if  $\Phi = \emptyset$  then
     $\Phi$  is satisfiable: return(satisfiable)
  if  $\mathbf{C}_k = \emptyset$  for a  $\mathbf{C}_k \in \Phi_{CNF}$  then
     $\Phi$  is contradictory: backtrack.
  if ( $\mathbf{Q}_k = \emptyset$  and  $\delta_k = -1$ ) for a  $\mathbf{Q}_k \in \Phi_{CoE}$  then
     $\Phi$  is contradictory: backtrack.
   $l := \text{branch\_rule}(\Phi)$ ;
  DPL( $\Phi \cup \{l\}$ , depth+1);
  DPL( $\Phi \cup \{\neg l\}$ , depth+1);
return(unsatisfiable)

```

Figure 1: The DPL algorithm extended for CNF/CoE formulas.

6. SOLVING THE DIMACS PARITY INSTANCES

We apply the techniques that we discussed previously to solve the DIMACS `par*-c.cnf` instances. These instances all contain a subformula with balanced polynomial representation. This subformula is a CNF translation of a CoE where all equivalency clauses have length three. We do not need to apply the LP approach to identify this formula, since the CoE subformula can be easily found by inspection. Note that if for example the order of the clauses would be changed, finding the CoE subformula in this way might no longer be practical; then the LP approach could be used.

All algorithms were implemented in C, and the results reported in this paper were obtained running the code on a SGI POWER CHALLENGE with a 200 Mhz R10k processor. All times reported are in seconds. In Table 1 we report on the results of the first phase of the algorithm which consists of isolating and solving the CoE subformulas. The initial numbers of variables and clauses are given by m and n . The number of equivalency clauses in the CoE subformula is denoted by k ; due to the specific structure of the instances, the number of clauses in the corresponding CNF equals $4k$, so the size of the remaining CNF is $n - 4k$ clauses. In the table we also indicate the number of independent variables determining the solutions of the CoE formula. The number of satisfying solutions for the CoE subformula equal $2^{|\mathcal{I}|}$. Note that the CoE formula does not need to be solved separately for the modified DPL algorithm to be valid. However, if it is solved, and subsequently it turns out that some dependent variable does not occur in the CNF part of the formula, this variable and the equivalency clause it occurs in need not be considered in the DPL search procedure. So, if we have the choice between two variables p_i and p_j of which only p_i occurs in the CNF subformula as well, we choose to remove

p_j from the set of independent variables. This allows us to reduce the problem size for phase two considerably.

Moreover, on solving the CoE formula an inconsistency might be detected. For example, the `dubois*.cnf` and `pret*.cnf` instances, which are also in the DIMACS suite, are already found to be unsatisfiable in the first phase of our algorithm. These instances are fully equivalent to CoE formulas and thus solved in polynomial time [9].

Before starting the second phase of the algorithm we first remove as many dependent variables and equivalency-clauses as possible. It may be noted that on branching strategies considering only the CNF subformula this has no effect as far as the node count is concerned; computation times however will reduce. The remaining numbers of variables, clauses and equivalency-clauses are given by m , n and k . Note that $m = k + |\mathcal{I}|$; each dependent variable occurs in exactly one equivalency-clause. We tested several branching strategies on the `par16*` instances, and used the one that appeared to be the best to solve the larger instances. In Table 2 we report on the results; all the instances are satisfiable. The branching strategy we arrived at is simply the *maximal occurrence in shortest clause* rule, with a lexicographic tie-break. We report on the node counts obtained by first branching to l and $\neg l$ respectively. The node count gives the number of times that a branching variable was chosen. A typical phenomenon of DPL algorithms that we also encountered here is that using different branching strategies the computation times and node counts may vary heavily.

Examining the tables we conclude that the smaller instances are solved in fractions of seconds, while the largest take at most about ten minutes.

The application of the techniques and notions described in this paper to more general SAT problems is the subject of further research.

instance	m	n	k	time	$ \mathcal{I} $
par8-1-c.cnf	64	254	56	.01	8
par8-2-c.cnf	68	270	60	.01	8
par8-3-c.cnf	75	298	67	.01	8
par8-4-c.cnf	67	266	59	.01	8
par8-5-c.cnf	75	298	67	.01	8
par16-1-c.cnf	317	1264	270	.08	47
par16-2-c.cnf	349	1392	302	.11	47
par16-3-c.cnf	334	1332	287	.09	47
par16-4-c.cnf	324	1292	277	.09	47
par16-5-c.cnf	341	1360	294	.10	47
par32-1-c.cnf	1315	5254	1158	5.68	157
par32-2-c.cnf	1303	5206	1146	5.00	157
par32-3-c.cnf	1325	5294	1168	5.76	157
par32-4-c.cnf	1333	5326	1176	5.61	157
par32-5-c.cnf	1339	5350	1182	5.80	157

Table 1: Results of the first phase of the algorithm

instance	m	n	k	nodes	time	nodes	time
par8-1-c.cnf	31	30	23	3	.00	1	.00
par8-2-c.cnf	31	30	23	3	.00	1	.00
par8-3-c.cnf	31	30	23	2	.00	1	.00
par8-4-c.cnf	31	30	23	3	.00	3	.00
par8-5-c.cnf	31	30	23	4	.00	4	.00
par16-1-c.cnf	124	184	77	82	.05	67	.04
par16-2-c.cnf	124	184	77	58	.03	144	.07
par16-3-c.cnf	124	184	77	55	.03	137	.07
par16-4-c.cnf	124	184	77	51	.03	131	.07
par16-5-c.cnf	124	184	77	49	.03	85	.05
par32-1-c.cnf	375	622	218	410634	492	130258	162
par32-2-c.cnf	375	622	218	201699	232	335988	418
par32-3-c.cnf	375	622	218	502747	615	6712	8
par32-4-c.cnf	375	622	218	218021	257	267032	347
par32-5-c.cnf	375	622	218	179325	221	328253	415

Table 2: Results of the second phase of the algorithm

References

1. S.A. Cook. The complexity of theorem proving procedures. In *Proceedings of the 3rd annual ACM symposium on the Theory of Computing*, pages 151–158, 1971.
2. J.M. Crawford, M.J. Kearns, and R.E. Schapire. The minimal disagreement parity problem as a hard satisfiability problem. Draft version, 1995.
3. M. Davis and H. Putnam. A computing procedure for quantification theory. *Journal of the ACM*, 7:210–215, 1960.
4. J. Gu. Global optimization for satisfiability (SAT) problem. *IEEE Transactions on Knowledge and Data Engineering*, 6(3):361–381, 1994.
5. M.A. Trick. Second DIMACS challenge test problems. In D.S. Johnson and M.A. Trick, editors, *DIMACS series in Discrete Mathematics and Theoretical Computer Science*, volume 26. American Mathematical Society, 1996.
6. N. Karmarkar. A new polynomial-time algorithm for linear programming. *Combinatorica*, 4:373–395, 1984.
7. D.W. Loveland. *Automated theorem proving: a logical basis*. North-Holland, Amsterdam, 1978.
8. B. Selman, H. Kautz, and D. McAllester. Ten challenges in propositional reasoning and search. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence (IJCAI-97)*, Nagoya, Aichi, Japan, 1997.
9. J.P. Warners and H. van Maaren. Satisfiability problems with balanced polynomial representation. Technical Report 97–47, Department of Technical Mathematics and Informatics, Faculty of Information Technology and Systems, Delft University of Technology, Delft, The Netherlands, 1997.