

The Complexity of Unique k -SAT: An Isolation Lemma for k -CNFs

Chris Calabro*
Russell Impagliazzo †
Valentine Kabanets ‡
Ramamohan Paturi §

June 20, 2006

Abstract

We provide some evidence that Unique k -SAT is as hard to solve as general k -SAT, where k -SAT denotes the satisfiability problem for k -CNFs with at most k literals in each clause and Unique k -SAT is the promise version where the given formula has 0 or 1 solutions. Namely, defining for each $k \geq 1$, $s_k = \inf\{\delta \geq 0 \mid \exists \text{ a } O(2^{\delta n})\text{-time randomized algorithm for } k\text{-SAT}\}$ and, similarly, $\sigma_k = \inf\{\delta \geq 0 \mid \exists \text{ a } O(2^{\delta n})\text{-time randomized algorithm for Unique } k\text{-SAT}\}$, we show that $\lim_{k \rightarrow \infty} s_k = \lim_{k \rightarrow \infty} \sigma_k$. As a corollary, we prove that, if Unique 3-SAT can be solved in time $2^{\epsilon n}$ for every $\epsilon > 0$, then so can k -SAT for all $k \geq 3$.

Our main technical result is an isolation lemma for k -CNFs, which shows that a given satisfiable k -CNF can be efficiently probabilistically reduced to a uniquely satisfiable k -CNF, with non-trivial, albeit *exponentially small*, success probability.

1 Introduction

While NP-complete search problems are presumably intractable to solve, this does not mean that all instances of these problems are difficult. Many heuristics for NP-complete problems, e.g., using back-tracking or local search techniques, take advantage of some structure in the instance. This leads to the general question: what makes an instance of a search problem easier than the worst-case? Conversely, we can ask, what types of structure cannot be exploited by heuristics? Here, we can consider both the structure of the instance (e.g. density of constraints to variables, or expansion of a formula when viewed as a hypergraph) and the structure of the solution space (e.g. how correlated are different solutions?).

One particularly natural measure of the solution space is the number of solutions. One can have different intuitions concerning how the number of solutions should affect the difficulty of the instance. On the one hand, one might suspect that the hardest instances would be those with few solutions, since they become “needles in a large haystack”. On the other hand, algorithms that use local search or narrow in on the solution

*University of California, San Diego

†University of California, San Diego. Research supported by NSF Award CCR-0098197 and USA-Israel BSF Grant 97-00188

‡Simon Fraser University, Canada. This research was done while the author was at the University of California, San Diego, supported by a postdoctoral fellowship from the Natural Sciences and Engineering Research Council of Canada

§University of California, San Diego. Research supported by NSF Award CCR-0098197

through deductions might do worse when there are many unrelated solutions. These algorithms might get caught between two equally attractive solutions, like Buridan’s ass (who starves to death, unable to decide between two sources of food).

The best known algorithm for general CNF satisfiability (without any constraints on the clause size) is little better than exhaustive search, running in time $O(2^{n(1-\frac{1}{lg\ m})})$, where m is the number of clauses [Schuler03]. The situation is slightly better for k -SAT, where each clause of a given CNF contains at most k literals, for some constant $k \geq 3$. This syntactic restriction can be exploited to obtain a $O(2^{\epsilon_k n})$ -time algorithm for k -SAT, for some constant $0 < \epsilon_k < 1$ dependent on k . The first such algorithms were obtained in the 1980’s [Dan81, MS85]. Recently, a number of faster algorithms have been proposed [PPZ99, PPSZ98, Sch02, HSSW02, DGH⁺]. The best known randomized algorithm for 3-SAT is given in [HSSW02]. For $k \geq 4$, the best known randomized algorithm is due to [PPSZ98]; its time complexity on k -CNFs with n variables is $O(2^{(1-\mu_k/(k-1))n})$, where $\mu_k > 1$ is an increasing function of k that approaches $\pi^2/6 \approx 1.644$.

The techniques of [PPSZ98] and its predecessor [PPZ99] crucially rely upon the notion of *isolation* of a satisfying assignment. The degree of isolation of a satisfying assignment $x \in \{0, 1\}^n$ is the number of coordinates $1 \leq i \leq n$ such that flipping the value of x_i turns x into a falsifying assignment. Thus, a unique satisfying assignment has degree n of isolation, whereas any satisfying assignment of a tautology has degree 0 of isolation. The argument in [PPZ99] proceeds as follows. If a given k -CNF has “many” satisfying assignments, then we are likely to find one by random sampling. On the other hand, if there are “few” satisfying assignments, then there must be a satisfying assignment of “high” degree of isolation. However, sufficiently isolated satisfying assignments have “short” descriptions (via the Satisfiability Coding Lemma in [PPZ99]), and so these assignments can be found by enumerating all possible “short” descriptions, and checking if any such description corresponds to a satisfying assignment. Combining the two cases yields a randomized algorithm of running time $\text{poly}(n)2^{(1-1/k)n}$. (Later, [PPSZ98] improved this algorithm.)

A closer look at the analysis in [PPZ99] reveals that the more satisfying assignments a given k -CNF has, the faster we can find one of them. More precisely, as we prove later in the paper (see Section 4), if a k -CNF on n variables has s satisfying assignments, then one of them can be found with probability $\geq \frac{1}{2}$ in time $\text{poly}(n)(2^n/s)^{1-1/k}$. This result suggests that the hardest case of k -SAT may be *Unique k -SAT*, the promise problem whose positive instances are *uniquely* satisfiable k -CNFs, and whose negative instances are unsatisfiable k -CNFs.

In the rest of our paper, we provide further evidence towards the truth of this hypothesis.

1.1 Related work.

Before describing our results, we would like to discuss how the result of Valiant and Vazirani [VV86] on Unique Formula SAT is related to our problem of Unique k -SAT; here, Unique Formula SAT is the promise problem with uniquely satisfiable Boolean formulas (not necessarily in conjunctive normal form) as its positive instances, and unsatisfiable formulas as its negative instances.

Valiant and Vazirani [VV86] show that in some formal sense, uniqueness of the solution cannot be used to solve search problems quickly. They give a randomized polynomial-time reduction from Formula SAT to instances of Formula SAT with unique solutions. In particular, they use a random polynomial time algorithm to transform a formula F to F' such that F' is satisfiable if and only if F is and if F is satisfiable, then F' is uniquely satisfiable with probability at least $\Omega(1/n)$. This shows that a probabilistic polynomial-time algorithm for Unique Formula SAT would also give a similar algorithm for general Formula SAT. The same follows for most other NP-complete problems, such as k -SAT, by combining this reduction with a reduction from Unique SAT to the problem (which usually preserves uniqueness).

However, there is a large gap between not being solvable in polynomial-time and being the “hardest case” of a problem. For example, it is consistent with the Valiant-Vazirani result that Unique k -SAT be solvable in $O(2^{\sqrt{n}})$ time, but general k -SAT requires $2^{\Omega(n)}$ time, where n is the number of variables in the formula. This is because, when combined with the reduction from Formula SAT to k -SAT, the Valiant-Vazirani reduction squares the number of variables in the formula. Due to the variable squaring, this reduction is not very useful if the best algorithm for Unique k -SAT runs in time $\Omega(2^{\epsilon_k n})$, for some $0 < \epsilon_k < 1$, which is the common working hypothesis. We need to use a different argument to show that k -SAT is no harder than Unique k -SAT.

We stress that our emphasis in this paper is on the *promise* problem Unique k -SAT, as opposed to the problem of *deciding* whether a k -CNF formula has a unique solution, as in [GB97]; the latter paper shows that this uniqueness decision problem for k -SAT is solvable in time $O(2^{\delta n})$ iff k -SAT is solvable in time $O(2^{\delta n})$.

1.2 Our Results

In the present paper, we considerably narrow the gap left by the Valiant-Vazirani reduction. We give a family of randomized reductions from general k -SAT to Unique k -SAT which succeed in reducing a satisfiable input to an instance of Unique k -SAT with probability at least $\Omega(2^{-\epsilon_k n})$ where $\epsilon_k \rightarrow 0$ as $k \rightarrow \infty$. This implies that if general k -SAT requires time $2^{\delta n}$ for some k and some $\delta > 0$, then, for any $\delta' < \delta$, Unique k' -SAT requires time $2^{\delta' n}$ for sufficiently large k' . Combining with results from [IPZ98], we show that if Unique 3-SAT does not require exponential time, then no search problem in SNP requires exponential time.

Before stating our main result, we need to state some key definitions that we use throughout the paper. Recall that k -SAT denotes the satisfiability decision problem for k -CNFs with at most k literals. Unique k -SAT denotes the promise decision problem of k -SAT where the number of solutions is either 0 or 1. Our results in this paper concern the asymptotic complexities of these problems as k increases. While we couch our results in terms of decision problems, they can naturally be extended to search versions of the same problems.

We will often speak of a randomized algorithm A ‘solving’ a problem L in time t . The exact meaning of this will be context dependent. If L is a *promise* problem (e.g., $L = \text{Unique-}k\text{-SAT}$), then we say that L is solved in randomized time t if there exists a randomized time t Turing machine returning ‘yes’ on positive input instances of L with probability at least $1/2$, and returning ‘no’ on negative input instances of L with probability 1.

When we talk about reductions, we mean Turing reductions. Thus we say that a (promise) problem L_1 (e.g., $L_1 = k\text{-SAT}$) is reducible to a (promise) problem L_2 (e.g., $L_2 = \text{Unique-}k\text{-SAT}$) in randomized time t if there is a randomized time t oracle Turing machine that solves L_1 when give oracle for L_2 (where, in case L_2 is a promise problem, the oracle for L_2 says ‘yes’ on positive instances of L_2 , and ‘no’ on negative instances of L_2 , but may say anything on a query outside the promise).

In this paper we will usually only care about the value of t up to a subexponential factor, i.e., up to a factor of $2^{\epsilon n}$ where $\epsilon > 0$ is a parameter, and so it is obviously sufficient to provide an algorithm running in time t with subexponentially small success probability, i.e., $\geq 2^{-\epsilon n}$, and then use repetition to boost the probability of success to be at least $1/2$.

We are solely concerned with the exponential complexity of k -SAT or Unique k -SAT. As our input complexity measure, we use n , the number of variables. We believe that the number of variables is a natural parameter as it measures the size of the space of solutions. All recent algorithms for k -SAT have running times expressed in the form $2^{\delta n}$, where $\delta > 0$ as a function of k approaches 1 as $k \rightarrow \infty$. Furthermore, the Sparsification Lemma proved in [IPZ98] (and stated explicitly in Lemma 6) shows that an arbitrary k -CNF

can be reduced in subexponential time to a union of at most subexponentially many k -CNFs each with a linear number of clauses. For a more thorough discussion of input complexity parameters, we refer the reader to [IPZ98].

With these considerations, our focus in this paper is on the quantities

$$s_k = \inf\{\delta \geq 0 \mid \exists \text{ a } O(2^{\delta n})\text{-time randomized algorithm for } k\text{-SAT}\}, \text{ and}$$

$$\sigma_k = \inf\{\delta \geq 0 \mid \exists \text{ a } O(2^{\delta n})\text{-time randomized algorithm for Unique } k\text{-SAT}\}$$

as functions of k . Our results show that the sequences s_k and σ_k have the same limit as $k \rightarrow \infty$. While our theorems hold unconditionally, they are more interesting in the case $\forall k \geq 3, s_k > 0$, which we call the Exponential Time Hypothesis (ETH).

Denoting $s_\infty = \lim_{k \rightarrow \infty} s_k$ and $\sigma_\infty = \lim_{k \rightarrow \infty} \sigma_k$, we can state our main result as follows.

Theorem 1. $s_\infty = \sigma_\infty$.

This gives strong evidence that Unique k -SAT is as hard to solve as general k -SAT. Combining Theorem 1 with the results in [IPZ98] yields the following.

Corollary 2. *If $\sigma_3 = 0$, then $s_\infty = 0$.*

The main technical ingredient in the proof of Theorem 1 is a more refined analysis of the random hashing technique of [VV86]. We modify the construction of [VV86] so that it uses *sparse* hyperplanes that depend on at most k' variables, where k' depends only on k and an approximation parameter ϵ . We then show that, given a satisfiable k -CNF ϕ , the modified construction is likely to produce a new k' -CNF ψ such that all satisfying assignments of ψ are clustered in a single Hamming ball of “small” radius. Finally, we argue that a random assignment of “few” randomly chosen variables of ψ will result in a uniquely satisfiable k' -CNF ψ' . This result is formally stated as Lemma 3 (Isolation Lemma) below.

It is worth pointing out that, unlike the reduction of [VV86], our randomized reduction has only *exponentially small* success probability. However, since we are interested in $O(2^{\epsilon n})$ -time algorithms, such small success probabilities suffice for our purposes.

Remainder of this paper In Section 2, we state and prove our main technical result, an Isolation Lemma for k -CNFs. Section 3 contains the proofs of Theorem 1 and Corollary 2 stated in the Introduction. In Section 4, we show that Unique k -SAT is the worst case for the k -SAT algorithm of [PPZ99].

2 Isolation Lemma for k -CNFs

Let $h(\epsilon) = -\epsilon \lg \epsilon - (1 - \epsilon) \lg(1 - \epsilon)$ denote the binary entropy function.

Lemma 3 (Isolation Lemma). *For every $k \in \mathbb{N}, \epsilon \in (0, \frac{1}{4})$, there is a randomized polynomial-time algorithm $I_{k,\epsilon}$ that, given an n -variable k -CNF ϕ , outputs an n -variable k' -CNF ϕ' , for $k' = \max\{\lceil \frac{1}{\epsilon} \ln \frac{2}{\epsilon} \rceil, k\}$, such that*

1. *if ϕ is unsatisfiable, then so is ϕ' , and*
2. *if ϕ is satisfiable, then, with probability at least $(32n2^{3h(2\epsilon)^n})^{-1}$, ϕ' has a unique satisfying assignment.*

Proof. Our proof will consist of two steps. First, we intersect the set S of satisfying assignments of ϕ with randomly chosen hyperplanes on k' variables, and show that, with reasonable probability, all the surviving satisfying assignments of ϕ are contained (concentrated) in a Hamming ball of “small” radius. Then we isolate a single satisfying assignment in this Hamming ball by randomly assigning a small number of coordinates.

STEP I: CONCENTRATION. Guess an integer $s \in [1..n]$, then with probability at least $1/n$, $\lg |S| \leq s \leq 1 + \lg |S|$. For the remainder of our argument, let us assume that our guess is correct. For m to be determined later, we intersect the solution space S of ϕ with m sparse hyperplanes of the Boolean cube $\{0, 1\}^n$, where each hyperplane h_i is chosen independently by the following randomized algorithm: Pick a subset $R_i \subseteq \{1, \dots, n\}$ of size k' uniformly at random. For each $j \in R_i$, pick $c_{i,j} \in \{0, 1\}$ uniformly at random. Pick $b_i \in \{0, 1\}$ uniformly at random. Output the hyperplane $h_i \stackrel{\text{def}}{=} \sum_{j \in R_i} c_{i,j} x_j = b_i$, where the summation is over the field \mathbb{F}_2 .

Let us write the system of linear equations h_1, \dots, h_m in the matrix form $Ax = b$, where A is an $m \times n$ 0-1 matrix, $x = (x_1, \dots, x_n)^t$, and $b = (b_1, \dots, b_m)^t$.

While a random sparse hyperplane is unlikely to separate two vectors $x, y \in \{0, 1\}^n$ that are close to each other, it has a reasonable chance of separating x and y that are Hamming distance $\text{dist}(x, y) \geq r \stackrel{\text{def}}{=} \lceil \epsilon n \rceil$ apart.

Claim 4. For any $x, y \in \{0, 1\}^n$ with $\text{dist}(x, y) \geq r$, we have

$$\Pr_A[Ax = Ay] \leq \frac{1}{2^m} e^{\epsilon m}$$

Proof of Claim 4. For any $x, y \in \{0, 1\}^n$ with $\text{dist}(x, y) \geq r$ and any $i \in \{1, \dots, m\}$, we have

$$\begin{aligned} \Pr[(Ax)_i \neq (Ay)_i] &= \Pr[(A(x - y))_i \neq 0] \\ &= \Pr[(A(x - y))_i \neq 0 \mid (x - y)_{R_i} \neq 0] \cdot \Pr[(x - y)_{R_i} \neq 0] \end{aligned}$$

where $(x - y)_{R_i}$ is the vector restricted to the coordinates from R_i . If $(x - y)_{R_i} \neq 0$, then $\Pr[(A(x - y))_i \neq 0] = \frac{1}{2}$. We get

$$\begin{aligned} \Pr[(Ax)_i \neq (Ay)_i] &= \frac{1}{2} \Pr[(x - y)_{R_i} \neq 0] \\ &= \frac{1}{2} \Pr[\exists j \in R_i (x - y)_j \neq 0] \\ &\geq \frac{1}{2} \left(1 - \frac{\binom{n-r}{k'}}{\binom{n}{k'}} \right) \\ &\geq \frac{1 - (1 - \epsilon)^{k'}}{2} \\ &\geq \frac{1 - (1 - \epsilon)^{(1/\epsilon) \ln(2/\epsilon)}}{2} \\ &\geq \frac{1 - e^{-\epsilon(1/\epsilon) \ln(2/\epsilon)}}{2} \\ &\geq \frac{1}{2} - \frac{\epsilon}{2}. \end{aligned}$$

Thus, $\Pr_A[Ax = Ay] \leq \left(\frac{1}{2} + \frac{\epsilon}{2}\right)^m \leq \frac{1}{2^m} e^{\epsilon m}$, as required. \square

Now, let us define the set of *semi-isolated* solutions as

$$S' = \{x \in S \mid \forall y \in S (\text{dist}(x, y) \geq r \rightarrow Ax \neq Ay)\}.$$

For each $x \in S$, we have

$$\begin{aligned} \Pr_A[x \in S'] &= 1 - \Pr_A[\exists y \in S \\ &\quad (\text{dist}(x, y) \geq r \wedge Ax = Ay)] \\ &\geq 1 - \sum_{y \in S, \text{dist}(x, y) \geq r} \Pr_A[Ax = Ay] \\ &\geq 1 - |S|e^{\epsilon m} / 2^m \end{aligned}$$

where the last inequality is by Claim 4. By setting $m = s + \lceil 4\epsilon s \rceil + 2$, we obtain

$$\Pr_A[x \in S'] \geq 1/2. \quad (1)$$

It follows that

$$\begin{aligned} &\sum_{x \in S} \Pr_{A,b}[x \in S' \wedge Ax = b] \\ &= \sum_{x \in S} \Pr_A[x \in S'] \Pr_b[Ax = b \mid x \in S'] \\ &\geq \frac{1}{2} |S| 2^{-m}, \end{aligned} \quad (2)$$

where the last inequality is by (1) and the fact that the vector b is chosen independently of the matrix A .

On the other hand, for any given A and b , the number of semi-isolated solutions $x \in S$ such that $Ax = b$ is at most $2^{h(\epsilon)n}$, since every pair of such solutions must be Hamming distance less than r apart. This implies

$$\sum_{x \in S} \Pr_{A,b}[x \in S' \wedge Ax = b] \leq 2^{h(\epsilon)n} \Pr_{A,b}[\exists x \in S (x \in S' \wedge Ax = b)]. \quad (3)$$

Combining inequalities (2) and (3) yields

$$\begin{aligned} \Pr_{A,b}[\exists x \in S' Ax = b] &\geq 2^{s-m-h(\epsilon)n-2} \\ &\geq 2^{-(4\epsilon+h(\epsilon))n-5}, \end{aligned} \quad (4)$$

concluding Step I of the lemma.

STEP II: ISOLATION. Now suppose that $\phi \wedge Ax = b$ is satisfiable and its solutions are in a Hamming ball of radius less than r , and hence, the diameter $d \leq 2\lceil \epsilon n \rceil$.

For assignments u and v to disjoint sets of variables, let uv denote the union of these assignments. Let uv and $u\bar{v}$ be two distinct solutions to $\phi \wedge Ax = b$ that are furthest apart. Let C be the set of variables on which uv and $u\bar{v}$ differ. Note that $|C| \leq d$. Fix $D \supseteq C$ such that $|D| = d$. Consider the assignment β_D to the variables in D consistent with the assignment uv : $x_i = v_i$ for $i \in C$ and $x_i = u_i$ for $i \in D - C$. The reduction chooses a random set D' of d variables and a random assignment $\alpha_{D'}$ to these variables. Define ϕ' to be the formula $\phi \wedge Ax = b \wedge x_{D'} = \alpha_{D'}$. If $D' = D$ and $\alpha_{D'} = \beta_D$, then uv is the *unique* solution to ϕ' , since any other solution wv to ϕ' would also be a solution to $\phi \wedge Ax = b$, but $\text{dist}(wv, u\bar{v}) > \text{dist}(uv, u\bar{v})$.

The probability of correctly guessing a set of variables D containing C and correctly assigning them is at least

$$\frac{2^{-d}}{\binom{n}{d}} \geq 2^{-(2\epsilon+h(2\epsilon))n}. \quad (5)$$

Combining equations (4) and (5), we conclude that the probability of correctly guessing s , A , b , D' , and $\alpha_{D'}$ so that the resulting formula $\phi' \stackrel{def}{=} \phi \wedge Ax = b \wedge x_{D'} = \alpha_{D'}$ is uniquely satisfiable is at least $\frac{1}{32n} 2^{-3h(2\epsilon)n}$, as claimed. \square

3 Applications of the Isolation Lemma

3.1 Proof of Theorem 1

Lemma 5. $s_k \leq \sigma_k + O(\frac{\ln^2 k}{k})$.

Proof. Let $\epsilon = \frac{2 \ln k}{k}$. Then $k \geq \frac{1}{\epsilon} \ln \frac{2}{\epsilon}$. For every $\gamma > 0$ and every k , we have a randomized algorithm $U_{k,\gamma}$ for Unique k -SAT that runs in time $O(2^{(\sigma_k+\gamma)n})$. Applying the Isolation Lemma 3, there is a polynomial time algorithm $I_{k,\epsilon}$ for reducing a given satisfiable k -CNF ϕ to a uniquely satisfiable k -CNF ϕ' with success probability $p \geq 2^{-O(h(\epsilon)n)} \geq 2^{-O(\frac{\ln^2 k}{k}n)}$. By running the algorithm $I_{k,\epsilon}$ $O(p^{-1})$ times and applying algorithm $U_{k,\gamma}$ to every output k -CNF ϕ' , we obtain a randomized algorithm for k -SAT that has running time $O(2^{(\sigma_k+\gamma+O(\frac{\ln^2 k}{k}))n})$. As $\gamma \rightarrow 0$, we have the lemma. \square

Theorem 1 then follows by taking the limit of the inequality in the previous lemma as k tends to infinity.

3.2 Proof of Corollary 2

We shall need the following result.

Lemma 6 (Sparsification Lemma [IPZ98]). *For all k and $\epsilon > 0$, there is a $\text{poly}(n)2^{\epsilon n}$ -time algorithm $S_{k,\epsilon}$ that, given an n -variable k -CNF ϕ , returns an n -variable k -CNF ψ that has at most $c(k, \epsilon)n$ clauses, for some function c of k, ϵ , and such that*

1. every solution to ψ is also a solution to ϕ , and
2. if ϕ is satisfiable, then so is ψ , with probability at least $\geq 2^{-\epsilon n}$.

We also need a simple clause-width-reducing algorithm R_3 which takes as input a k -CNF ϕ in n variables and m clauses and returns a 3-CNF ψ in $n + (k - 3)m$ variables whose solution set has a trivial bijective correspondence to that of ϕ . We describe R_3 . For each clause $l_1 \vee \dots \vee l_k$, where the l_i are literals, introduce $k - 3$ new variables y_3, \dots, y_{k-1} . The idea is to express $l_1 \vee \dots \vee l_k$ as

$$(l_1 \vee l_2 \vee y_3) \wedge (y_3 \vee \bar{l}_3) \wedge \dots \wedge (y_3 \vee \bar{l}_k) \wedge (\bar{y}_3 \vee l_3 \vee \dots \vee l_k),$$

and then recurse on the last clause. The above is equivalent to

$$(l_1 \vee l_2 \vee y_3) \wedge (y_3 \leftrightarrow (l_3 \vee \dots \vee l_k)),$$

and so we can see that for each solution x of ϕ , there is exactly one solution y to ψ and that we can compute x from y simply by throwing away the new variables.

Proof of Corollary 2. We combine sparsification, isolation, and clause-width reduction. Suppose that for each $\epsilon > 0$, there is a randomized algorithm $U_{3,\epsilon}$ that solves Unique 3-SAT in time $O(2^{\epsilon n})$.

Let $k \geq 3$ be arbitrary. We want to show that $s_k = 0$. For any $\epsilon \in (0, \frac{1}{4})$, the algorithm $T \stackrel{\text{def}}{=} R_3 \circ I_{k,\epsilon} \circ S_{k,\epsilon}$, composed of the three algorithms R_3 , $I_{k,\epsilon}$, and $S_{k,\epsilon}$, returns a 3-CNF with at most $c(k, \epsilon)n$ variables for some function c of k, ϵ since the construction in the Isolation Lemma adds at most $c'(k, \epsilon)n$ clauses to the sparsified formula, where $c'(k, \epsilon)$ is a function of k and ϵ . Composing $U_{3,\epsilon}$ with T , we obtain a $O(2^{2\epsilon n})$ -time algorithm for k -SAT that has success probability $p \geq 2^{-O(h(\epsilon)n)}$. Running this algorithm $O(p^{-1})$ times independently, we get a randomized algorithm for k -SAT of running time $2^{O(h(\epsilon)n)}$. Letting $\epsilon \rightarrow 0$ gives the result. \square

4 The More Satisfying Assignments, the Better

Here we present a randomized algorithm for k -SAT whose expected running time on an n -variable poly(n)-size k -CNF with s satisfying assignments is $\text{poly}(n)(2^n/s)^{1-1/k}$. The algorithm is the same as that in [PPZ99]; our analysis of this algorithm is somewhat more refined.

Below, we say that a variable x is *forced* in a CNF ϕ if ϕ contains a unit clause which consists of x or \bar{x} only. Clearly, the forced variable will have the same value in every satisfying assignment to ϕ .

Consider the algorithm **A** given in Figure 4. Clearly, algorithm **A** is polynomial-time. We shall lower-

Algorithm A

INPUT: k -CNF ϕ on variables x_1, \dots, x_n

1. Pick a permutation σ of the set $\{1, \dots, n\}$ uniformly at random.
2. For $i = 1$ to n , if the variable $x_{\sigma(i)}$ is forced, then set $x_{\sigma(i)}$ so that the corresponding unit clause is satisfied; otherwise, set $x_{\sigma(i)}$ to 1 (TRUE) or 0 (FALSE) uniformly at random. Delete satisfied clauses and false literals.
3. If the constructed assignment a satisfies ϕ , then output a .

bound the probability that this algorithm succeeds at finding a satisfying assignment.

Lemma 7. *Let ϕ be a k -CNF on n variables that has s satisfying assignments. Then*

$$\Pr[\mathbf{A}(\phi) \text{ outputs some satisfying assignment}] \geq (s2^{-n})^{1-1/k}.$$

Proof. Let S be the set of satisfying assignments of ϕ , $a \in S$, and let σ be any permutation. For this fixed σ , consider the run of **A** that results in outputting the assignment a . We denote by $d_\sigma(a)$ the number of variables that were *not* forced during this run of **A**. Since the truth value of each non-forced variable is chosen uniformly at random, we have $\Pr[\mathbf{A}(\phi) = a \mid \sigma] = 2^{-d_\sigma(a)}$, where the probability is over the random choices of the algorithm **A** for non-forced variables, given the fixed permutation σ .

It is easy to see that

$$\begin{aligned} & \Pr[\mathbf{A}(\phi) \text{ outputs some satisfying assignment}] \\ &= \sum_{a \in S} \mathbf{E}_\sigma[2^{-d_\sigma(a)}] \geq \sum_{a \in S} 2^{-\mathbf{E}_\sigma[d_\sigma(a)]}, \end{aligned}$$

where the last step uses Jensen's inequality.

We can bound $\mathbf{E}_\sigma[d_\sigma(a)]$ using the argument from [PPZ99]. Let $l(a)$ denote the number of satisfying assignments a' of ϕ such that the Hamming distance $\text{dist}(a, a') = 1$. Then there are $n - l(a)$ variables such that each of them occurs as a unique true literal in some clause; such a clause is called *critical* in [PPZ99]. It follows that each such variable $x_{\sigma(i)}$ will be forced for a random permutation σ if $x_{\sigma(i)}$ occurs last in the corresponding critical clause. The latter happens with probability at least $1/k$ since each clause of ϕ is of size at most k . By the linearity of expectation, we obtain that the expected number of *forced* variables is at least $(n - l(a))/k$. Hence, $\mathbf{E}_\sigma[d_\sigma(a)] \leq n - (n - l(a))/k$.

Thus,

$$\begin{aligned} & \Pr[\mathbf{A}(\phi) \text{ outputs some satisfying assignment}] \\ & \geq \sum_{a \in S} 2^{-n(1-1/k)} 2^{-l(a)/k}. \end{aligned}$$

After factoring out $2^{-n(1-1/k)}$, we need to lower-bound the sum $\sum_{a \in S} 2^{-l(a)/k}$. Let $L = \mathbf{E}_{a \in S}[l(a)]$. By Jensen's inequality, we obtain

$$\sum_{a \in S} 2^{-l(a)/k} \geq s 2^{-L/k}, \quad (6)$$

where $s = |S|$.

Finally, to bound L , we use the well-known edge isoperimetric inequality due to [Har67]. This inequality implies that, for any subset S' of the n -dimensional Boolean cube $\{0, 1\}^n$,

$$|\{(a, a') \mid a, a' \in S' \text{ and } \text{dist}(a, a') = 1\}| \leq |S'| \lg |S'|.$$

Hence, for $S' = S$, we obtain $\sum_{a \in S} l(a) \leq s \lg s$, and so, $L = \mathbf{E}_{a \in S}[l(a)] \leq \lg s$.

Continuing inequality (6), we get

$$\sum_{a \in S} 2^{-l(a)/k} \geq s 2^{-(\lg s)/k} = s^{1-1/k}.$$

So, $\Pr[\mathbf{A}(\phi) \text{ outputs some satisfying assignment}] \geq (s 2^{-n})^{1-1/k}$, as required. \square

Lemma 7 immediately implies the following.

Theorem 8. *There is a randomized algorithm for k -SAT that, given a satisfiable k -CNF ϕ on n variables with s satisfying assignments, will output one of the satisfying assignments of ϕ in expected time $\text{poly}(n)(2^n/s)^{1-1/k}$.*

5 Open Problems

We would like to argue that the true challenge in designing an improved k -SAT algorithm is to solve *Unique k -SAT* faster than the corresponding algorithm of [PPSZ98].

While we are able to show that $s_\infty = \sigma_\infty$, we conjecture that $\forall k, s_k = \sigma_k$, which would be strong evidence that Unique k -SAT should be our focus. However, our techniques do not lend themselves to proving this conjecture. To prove the conjecture, it would be sufficient to prove an Isolation Lemma for k -CNF with an inverse polynomial probability of success. While it takes a family of size $\Omega(2^s)$ to isolate an element of an arbitrary set of size 2^s by intersecting with a randomly chosen set from the family, it would be

interesting if one can construct polynomial size families to isolate elements from sets of satisfying solutions of k -CNF.

In fact, we could hope to prove an even stronger version of the conjecture to the effect that satisfiability for formulas with many satisfying assignments is strictly easier than for the formulas with fewer solutions. To prove such a result, one might show that adding additional random constraints to a formula with many satisfying assignments preserves satisfiability while implicitly reducing the number of variables. For example, in our proof of the Isolation Lemma, we added a number of random sparse linear equations, which implicitly define some variables in terms of the others. Can we show that these constraints can be used to simplify search? A less ambitious formulation would be to prove the above just for Davis-Putnam algorithms and their variants.

References

- [Coo71] S.A. Cook. The complexity of theorem-proving procedures. In *Proceedings of the Third Annual ACM Symposium on Theory of Computing*, pages 151–158, 1971.
- [Dan81] E. Dantsin. Two propositional proof systems based on the splitting method. *Zapiski Nauchnykh Seminarov LOMI*, 105:24–44, 1981. (in Russian). English translation in *Journal of Soviet Mathematics*, 22(3):1293–1305, 1983.
- [DGH⁺] E. Dantsin, A. Goerdt, E.A. Hirsch, R. Kannan, J. Kleinberg, C. Papadimitriou, P. Raghavan, and U. Schöning. A deterministic $(2 - 2/(k + 1))^n$ algorithm for k -SAT based on local search. *Theoretical Computer Science*. to appear.
- [GB97] E. Grandjean, H. Buning. SAT-problems and reductions with respect to the number of variables. *Journal of Logic and Computation*, 7(4):457–471, 1997.
- [Har67] L.H. Harper. A necessary condition on minimal cube numberings. *Journal of Applied Probability*, 4:397–401, 1967.
- [HSSW02] T. Hofmeister, U. Schöning, R. Schuler, and O. Watanabe. A probabilistic 3-SAT algorithm further improved. In *Proceedings of the Nineteenth Annual Symposium on Theoretical Aspects of Computer Science*, pages 192–202, 2002.
- [IP01] R. Impagliazzo, R. Paturi. On the complexity of k -SAT. *Journal of Computer and Systems Sciences*, 62(2):367–375, 2001.
- [IPZ98] R. Impagliazzo, R. Paturi, and F. Zane. Which problems have strongly exponential complexity? In *Proceedings of the Thirty-Ninth Annual IEEE Symposium on Foundations of Computer Science*, pages 653–662, 1998.
- [MS85] B. Monien and E. Speckenmeyer. Solving satisfiability in less than 2^n steps. *Discrete Applied Mathematics*, 10:287–295, 1985.
- [PPSZ98] R. Paturi, P. Pudlák, M.E. Saks, and F. Zane. An improved exponential-time algorithm for k -SAT. In *Proceedings of the Thirty-Ninth Annual IEEE Symposium on Foundations of Computer Science*, pages 628–637, 1998.

- [PPZ99] R. Paturi, P. Pudlák, and F. Zane. Satisfiability coding lemma. *Chicago Journal of Theoretical Computer Science*, 1999. (preliminary version in FOCS'97).
- [Sch02] U. Schöning. A probabilistic algorithm for k -SAT based on limited local search and restart. *Algorithmica*, 32:615–623, 2002.
- [Schuler03] R. Schuler. An algorithm for the satisfiability problem of formulas in conjunctive normal form. Technical Report, Universit at Ulm, 2003.
- [VV86] L. Valiant and V. Vazirani. NP is as easy as detecting unique solutions. *Theoretical Computer Science*, 47:85–93, 1986.