# ANALYSIS OF TWO SIMPLE HEURISTICS

# ON A RANDOM INSTANCE OF k−SAT

by

Alan Frieze
Department of Mathematics
Carnegie Mellon University
Pittsburgh, PA 15213, U.S.A.

and

Stephen Suen
Department of Mathematics
Carnegie Mellon University
Pittsburgh, PA 15213, U.S.A.

# ANALYSIS OF TWO SIMPLE HEURISTICS ON A RANDOM INSTANCE OF $k$-SAT

Alan Frieze*and Stephen Suen
Department of Mathematics,
Carnegie Mellon University,
Pittsburgh PA15213, U.S.A.

December 5, 1992

## Abstract

We consider the performance of two algorithms, GUC and SC studied by Chao and Franco [2], [3], and Chvátal and Reed [4], when applied to a random instance $\omega$ of a boolean formula in conjunctive normal form with $n$ variables and $\lfloor cn \rfloor$ clauses of size $k$ each. For the case where $k = 3$, we obtain the exact limiting probability that GUC succeeds. We also consider the situation when GUC is allowed to have limited backtracking, and we improve an existing threshold for $c$ below which almost all $\omega$ is satisfiable. For $k \geq 4$, we obtain a similar result regarding SC with limited backtracking.

# 1   Introduction

Given a boolean formula $\omega$ in conjunctive normal form, the *satisfiability problem* (SAT) is to determine whether there is a truth assignment that satisfies $\omega$. Since SAT is NP-complete, one is interested in efficient heuristics that perform well "on average," or with high probability. The choice of the probabilistic space is crucial for the significance of such a study. In particular, it is easy to decide SAT in probabilistic spaces that generate formulas with large clauses [7]. To circumvent this problem, recent studies have focused on formulas with exactly $k$ literals per clause (the $k$-SAT problem). Of particular interest is the case $k = 3$, since this is the minimal $k$ for which the problem is NP-complete.

Let $V_n$ be a set of $n$ variables. Consider the space $\Omega_{m,n}^{(k)}$ of all $m = \lfloor cn \rfloor$ clause formulas over the variables in $V_n$ with exactly $k$ literals per clause. We assume that each clause uses

1

$k$ distinct variables to avoid trivial complications. Experimental evidence [8, 9] strongly suggests that there exists a threshold $\gamma$, such that formulas are almost surely satisfiable for $c < \gamma$ and almost surely unsatisfiable for $c > \gamma$, where $\gamma$ is about 4.2. This has not been proven rigorously, but such a threshold (namely c=1) is known to exist for 2-CNF formulas [6, 4].

Most practical algorithms for the satisfiability problem (such as the well-known Davis-Putnam algorithm [5]) work iteratively. At each iteration, the algorithm selects a literal and assigns it the value 1. All clauses containing this literal are erased from the formula, and the complement of the chosen literal is erased from the remaining clauses. Algorithms differ in the way they select the literal for each iteration. The following three rules are the most common ones:

1. *The unit clause rule:* If a clause contains only one literal, that literal must have the value 1;

2. *The pure literal rule:* If a formula contains a literal but does not contain its complement, this literal is assigned the value 1;

3. *The smallest clause rule:* Give value 1 to a (random) literal in a (random) smallest clause.

Broder, Frieze and Upfal [1] analysed an algorithm based entirely on the pure literal rule. They showed that in the $\Omega_{m,n}^{(3)}$ probabilistic space, the pure literal rule alone is sufficient to find, with high probability, a satisfying assignment for a random formula $\omega \in \Omega_{m,n}^{(3)}$, for $c = m/n \leq 1.63$. On the other hand, if $c > 1.7$, then the pure literal rule by itself does not suffice.

Chao and Franco [2],[3] and Chvátal and Reed [4] analysed two heuristics GUC and SC based on the small clause rule:
**begin**
   **repeat**
      choose a literal $x$;
      remove all clauses from $\omega$ that contain $x$ and remove $\bar{x}$ from any remaining clause;
      if a clause becomes empty - HALT, FAILURE;
   **until** no clauses left;
   HALT, SUCCESS
**end**

The algorithms GUC and SC differ in how the literal $x$ is chosen. In GUC, $x$ is chosen at random from a randomly selected clause of smallest size. SC (see Chvátal and Reed [4] for a complete description of SC) differs from GUC in that if there are no clauses of size one or two, then $x$ is chosen at random (or arbitrarily) from the set of all free literals. Since at least one clause is satisfied each time when GUC assigns a value to a variable, it is intuitively clear that GUC is likely (probabilistically) to perform better than SC. Algorithm SC however has

the advantage of being simpler to analyse. The reason for this is that since SC only takes care of clauses of size one and two, there are fewer cases to consider when analysing SC.

The combined results (among other things) in Chao and Franco [2], [3] and Chvátal and Reed [4] can be summarized as follows. For 3-SAT, if $c < 2/3$ then $SC$ succeeds with probability tending to 1 [4] and if $c < 2.99$ then the probability that UC (a variant of GUC using only the unit clause rule) succeeds does not tend to zero [2]. For $k$-SAT where $k \geq 4$, if

$$c < \left(\frac{k-1}{k-3}\right)^{k-3} \frac{k-1}{k-2} \frac{2^{k-3}}{k},$$

then SC succeeds with probability tending to 1 [4], and if $k \leq 40$ and

$$c < 0.7725 \left(\frac{k-1}{k-2}\right)^{k-2} \frac{2^k}{k+1}, \tag{1}$$

then the probability that GUC succeeds does not tend to zero [3].

Our first theorem gives the *precise* limiting probability that GUC succeeds when applied to a random instance of 3-SAT. Let $c_3 \approx 3.003$ be the solution to the equation

$$3c - 2\log c = 6 - 2\log(2/3),$$

and

$$f(x) = f_c(x) = \frac{3c}{4}(1 - x^2) + \log x, \qquad x \in (0, 1).$$

When $c < c_3$ we have $f(x) < 1$ for all $x \in (0, 1)$.

**Theorem 1** *Consider applying GUC to a random instance of 3-SAT with $n$ variables and $\lfloor cn \rfloor$ clauses.*
*(a) Suppose that $c < 2/3$. Then*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{GUC succeeds}) = 1.$$

*(b) Suppose that $2/3 \leq c < c_3$. Let $\alpha$ be the unique root of $f(x) = 0$ that is strictly less than 1. Then*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{GUC succeeds}) = \exp\left(-\int_\alpha^1 \frac{f(x)^2}{4x(1 - f(x))} dx\right).$$

*(c) If $c \geq c_3$ then*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{GUC succeeds}) = 0.$$

Chao and Franco [2] report that using GUC in a backtracking algorithm can be quite successful (and possibly be polynomial expected time for certain values of $c$). We describe (in Section 5) a modification of GUC called GUCB that allows a limited amount of backtracking when an empty clause is produced. We obtain the following result by showing that for sufficiently small $c$, the backtracking does not change the state of GUC by a great deal.

3

**Theorem 2** *Consider GUCB when applied to a random instance of 3-SAT with $n$ variables and $\lfloor cn \rfloor$ clauses. If $c < c_3$ then*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{GUCB succeeds}) = 1.$$

We next turn our attention to algorithm SC. It is possible to show that the assertions in Theorems 1 and 2 hold for SC. In fact, our proof of Theorem 1 can be extended to obtain the precise limiting probability that SC succeeds when applied to a random instance of $k$-SAT. However, the more interesting question is: for what values of $c$ will SC, with limited backtracking as in GUCB, succeed with probability close to 1? We answer this question with our next result.

Assume $k \geq 4$. Let

$$p_3(x) = \binom{k}{3} \frac{c}{2^{k-3}} x^2 (1-x)^{k-3}.$$

It is easy to see that $p_3(x)$ is unimodal, achieving a maximum of

$$\frac{2}{3} \frac{kc}{2^{k-3}} \frac{k-2}{k-1} \left( \frac{k-3}{k-1} \right)^{k-3}$$

when $x = 2/(k-1)$. For

$$c > \frac{2^{k-3}}{k} \left( \frac{k-1}{k-3} \right)^{k-3} \frac{k-1}{k-2},$$

let $0 < \beta_0 = \beta_0(c) < \beta_1 = \beta_1(c) < 1$ be the two solutions of the equation $p_3(x) = 2/3$. We prove the following theorem.

**Theorem 3** *Suppose that $k \geq 4$. Let $c_k$ be the maximum value of $c$ such that*

$$\frac{1}{(k-1)(k-2)} \left( \frac{1}{\beta_0^2} + \frac{k-3}{\beta_0} - \frac{1}{\beta_1^2} - \frac{k-3}{\beta_1} \right) + \ln(\beta_0/\beta_1) \leq 1.$$

*Then when SCB is applied to a random instance of $k$-SAT with $n$ variables and $\lfloor cn \rfloor$ clauses where $c < c_k$, we have*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{SCB succeeds}) = 1.$$

Write $c_k = \eta_k 2^k / k$. It is possible to show that as $k \to \infty$, $\eta_k \to \eta^*$ where $\eta^*$ can be defined similarly as $\eta_k$. Numerical calculations show that $\eta^* \approx 1.817$, $\eta_4 \approx 1.3836$, $\eta_5 \approx 1.504$, $\eta_{10} \approx 1.686$, and that $\eta_k$ is increasing in $k$. Theorem 3 gives a constant $c_k$ such that almost every formula $\omega$ with $n$ variables and $\lfloor cn \rfloor$, with $c < c_k$, clauses of size $k$ is satisfiable. This improves, by only a constant factor, a similar result in [4]. Also, $c_k$ (for $4 \leq k \leq 40$) is smaller than the right hand side of (1), and we believe that if the limiting probability that GUC succeeds is positive, then GUC with limited backtracking (as described later) succeeds with probability $1 - o(1)$. It is thus very likely that when applied to random instances of

4

$k$-SAT for $k \geq 4$, GUCB has a higher threshold of success than SCB. At present, we can only characterize the critical behaviour of GUC and GUCB, when applied to random instances of $k$-SAT with $k \geq 4$, using a system of $k - 2$ polynomial equations whose properties we have difficulty in penetrating analytically. It seems unlikely that the exact thresholds for GUCB can be rid of the factor $1/k$ (see definition of $c_k$).

The basis of our proof of Theorem 1 is that the intermediate states of GUC (or SC), when applied to a random instance of $k$-SAT, can be represented by a Markov chain which we describe as follows. Consider GUC when applied to a formula $\omega$ chosen at random (with equal probability) from the space $\Omega_{m,n}^{(k)}$ where $m = \lfloor cn \rfloor$. For $V \subseteq V_n$, let $C_j(V)$ be the set of all clauses of size $j$ chosen from a set $V$ of variables. Use $\nu$ to denote the number of variables whose truth values are not yet determined by GUC at an intermediate stage. We call this stage $\nu$ and so GUC starts at stage $n$. For the purpose of analysis, all empty clauses are assumed to be removed by GUC as soon as they are created, and GUC is allowed to run until the set of clauses is exhausted. Hence, GUC succeeds if and only if the number of empty clauses created is zero. Let $N_i = N_i(\nu)$ $(i = 0, 1, 2, \ldots, k)$ be the number of clauses containing $i$ literals at stage $\nu$ of GUC. Thus $N_i$ is a function of $\nu$. Note that if $V_\nu$ is the set of variables whose truth values remain unassigned by GUC at stage $\nu$, then a clause with size $i$ $(1 \leq i \leq k)$ in the remaining set of clauses at stage $\nu$ is equally likely to be a clause in $C_k(V_\nu)$ (since each clause initially is equally likely to be any clause in $C_i(V_n)$). Hence $N = (N_0, N_1, \ldots, N_k)$ is a Markov chain.

We next write down the transition probability of $N$. Use $B(\tau, p)$ to denote a binomial variable with parameters $\tau$ and $p$ and note that $\nu$ decreases by 1 at each stage. Write $\Delta N_i(\nu) = N_i(\nu - 1) - N_i(\nu)$ as the change from stage $\nu$ to stage $\nu - 1$. Then $\Delta N_i$ are binomial variables (conditional upon $N(\nu)$). We shall write down the distributions of $\Delta N_i$ under the different cases where the minimum size of the clauses is $i$. For $i = 1, 2, \ldots, k$, we write $\chi_i((y_0, y_1, \ldots, y_k)) = 1$ if $\min\{j \mid y_j \neq 0, 1 \leq j \leq k\} = i$, and $\chi_i((y_0, y_1, y_2, \ldots, y_k)) = 0$ if otherwise. Also, $\chi_0(y) = 0$ always. Consider the stage $\nu$ when GUC has just assigned 1 to a literal $x$ in clause $C$ and is about to remove clauses that contain $x$ and all occurrences of $\bar{x}$ from other clauses. Let $\Delta_{j,0}$ be the number of clauses of size $j$ containing literals $x$ or $\bar{x}$ (but not including $C$). Let $\Delta_{j,1}$ be the number of clauses with size $j$ containing literal $\bar{x}$ (but not $x$ as all variables in a clause are different). It is simple to check that given $N = N(\nu)$, we have for $j = 1, 2, \ldots, k$ that

$$
\begin{aligned}
\Delta_{j,0}(\nu) &= B(N_j - \chi_j(N), j/\nu) \\
\Delta_{j,1}(\nu) &= B(\Delta_{j,0}, 1/2)
\end{aligned}
$$

Then for $j = 0, 1, \ldots, k$,

$$
\Delta N_j(\nu) = \Delta_{j+1,1}(\nu) - \Delta_{j,0}(\nu) - \chi_j(N(\nu)),
$$

where $\Delta_{0,0} = \Delta_{k+1,1} = 0$. Note that if $N_1(\nu) = 0$, then $\Delta N_0(\nu) = 0$ with probability 1. Note also that if $N_1(\nu) \geq 1$, then a clause of size one (with literal $x$ say) is chosen at stage $\nu$ and the probability that $\Delta N_0(\nu) = 0$ equals the probability that none of the other $N_1(\nu) - 1$

clauses contains the literal $\bar{x}$, which is precisely $(1 - 1/2\nu)^{N_1-1}$. Hence,

$$\mathbf{Pr}(\text{GUC succeeds}) = \mathbf{E}\left[\prod_{\nu=n_1}^{n}\left(1 - \frac{1}{2\nu}\right)^{(N_1(\nu)-1)^+}\right], \tag{2}$$

where $n_1$ is the stage when all clauses are removed. Theorem 1(b) is obtained by finding accurate estimates for $\sum_\nu(N_1(\nu) - 1)^+/(2\nu)$ in the case of 3-SAT. Theorem 1(a) and (c) are shown using monotonicity arguments.

We shall also require similar statements for SC. Let $N'_j(\nu)$ be the number of size $j$ clauses remaining at stage $\nu$ when SC is applied to a random instance of $k$-SAT with $n$ variables and $m$ clauses. Then similar to GUC, $N'(\nu)$ is a Markov chain with initial state $N'(n) = (0,\ldots,0,m)$ and transition probabilities given by

$$\Delta N'_j(\nu) = \begin{cases} \Delta'_{j+1,1}(\nu) - \Delta'_{j,0}(\nu) - \chi_j(N'(\nu)), & \text{if } j = 0,1,2, \\ \Delta'_{j+1,1}(\nu) - \Delta'_{j,0}(\nu), & \text{otherwise,} \end{cases}$$

where $\Delta'_{0,0} = \Delta'_{k+1,1} = 0$ and for $j = 1,2,\ldots,k$

$$\begin{aligned} \Delta'_{j,0}(\nu) &= \begin{cases} B(N'_j - \chi_j(N'), j/\nu), & \text{if } j = 0,1,2, \\ B(N'_j, j/\nu), & \text{otherwise,} \end{cases} \\ \Delta'_{j,1}(\nu) &= B(\Delta'_{j,0}, 1/2). \end{aligned}$$

Similar to (2), we have

$$\mathbf{Pr}(\text{SC succeeds}) = \mathbf{E}\left[\prod_{\nu=n'_1}^{n}\left(1 - \frac{1}{2\nu}\right)^{(N'_1(\nu)-1)^+}\right], \tag{3}$$

where $n'_1$ is the stage when all clauses are removed.

The layout of this paper is as follows. We concentrate on showing Theorems 1 and 2, while we shall only sketch our proof of Theorem 3. In the next section, we collect some useful properties of a Markov chain $X_t$ which will be used to approximate $N_1$ in proving Theorem 1(b). We shall then prove parts (a) and (c) of Theorem 1 in Section 3 by developing monotonicity arguments for comparing different Markov chains. Theorem 1(b) is proved in Section 4 by applying the results stated in Section 2. In Section 5, we describe how GUC is allowed to backtrack, and prove Theorem 2. In Section 6, we sketch briefly how our proof of Theorem 2 can be extended to proving Theorem 3.

## 2  A Markov chain

Use $B(m,p)$ to denote a binomial variable with parameters $m$ and $p$, and write $b_j = b_j(m,p)$ for the probability that $B(m,p)$ equals $j$. We assume throughout this section that $mp \leq \lambda^* < 1$. The big O terms in this section are uniform in $m$ and $p$ (but may depend on $\lambda^*$).

We consider a Markov chain $X_t$ with transition probabilities defined as follows. If $X_t = 0$, then $\Delta X_t = X_{t+1} - X_t$ equals $B(m,p)$ in distribution; otherwise $\Delta X_t$ equals $B(m,p) - 1$ in distribution. We assume $X_0 \geq 0$ and so $X = 0$ is a reflecting barrier. As we are interested in bounds that are uniform in $m$ and $p$, we need to consider a Markov chain $Y_t$ which is similar to $X_t$ except that in the one-step transitions of $Y_t$, we have a Poisson variable $P(\lambda)$ in place of $B(m,p)$. It will be clear that the two chains $X_t$ and $Y_t$ are very similar when $mp = \lambda$, although it is not possible to couple them so that $X_0 = Y_0$ and $X_t \leq Y_t$ for all $t \geq 0$. We let $\lambda = mp$ in this section.

Note that $X_t$ has a steady state distribution, denoted by $\pi$, satisfying

$$\pi_i = \pi_0 b_i + \sum_{j=1}^{i+1} \pi_j b_{i-j+1}, \quad \forall i \geq 0.$$

Writing $G_X(s) = \sum_{i=0}^{\infty} s^i \pi_i$ as the probability generating function of the steady state distribution, it follows from the above equations that

$$
\begin{aligned}
G_X(s) &= \pi_0 \sum_{i \geq 0} s^i b_i + \sum_{j \geq 1} \pi_j s^{j-1} \sum_{i \geq 0} b_i s^i \\
&= \pi_0 (1 - p + ps)^m + \frac{1}{s}(G_X(s) - \pi_0)(1 - p + ps)^m,
\end{aligned}
$$

giving

$$G_X(s) = \frac{\pi_0(s - 1)}{s(1 - p + ps)^{-m} - 1}.$$

As $G_X(1) = 1$, we have $\pi_0 = 1 - mp$ and

$$G_X(s) = \frac{(s-1)(1 - mp)}{s(1 - p + ps)^{-m} - 1}. \tag{4}$$

Since $(1 - p + ps)^m \leq \exp(-\lambda + \lambda s)$ for all $s$, we see that

$$G_X(s) \leq G(s) = \frac{(s-1)(1 - \lambda)}{s\exp(\lambda - \lambda s) - 1}, \tag{5}$$

for all $s$ between 1 and the radius of convergence of $G$. (It can be checked that $G(s)$ is the probability generating function of the steady state distribution of $Y_t$.) Since $\lambda \leq \lambda^*$, $G(s)$ exists for all $s < r_1^*$, where $r_1^* > 1$ is a constant depending on $\lambda^*$ only. ($r_1^*$ is in fact the unique root bigger than 1 of $s\exp(\lambda^* - \lambda^* s) = 1$.) Thus, (5) holds for all $s$ satisfying $1 < s < r_1^*$. Note also that from (4), the mean of the steady state distribution of $X_t$ is

$$\mu = \mu(m,p) = \sum_{i \geq 0} i\pi_i = \frac{mp(2 - p - mp)}{2(1 - mp)}. \tag{6}$$

We would like to consider the number of times that $X_t$ returns to 0 in a certain time period. To do this, we need to collect some preliminary results. Suppose $X_0 = 1$. Let $H_X$ be the time elapsed when $X_t$ first hits 0. ($H$ is defined accordingly for $Y_t$ with $Y_0 = 1$.) Note that

$$H_X = 1 + L_1 + \ldots + L_B \quad \text{in distribution},$$

where $B = B(m, p)$ in distribution and $L_1, \ldots, L_B$ are independent copies of $H_X$. Hence, writing $M_X(\theta) = \mathbf{E}[\exp(\theta H_X)]$, we have

$$M_X(\theta) = e^\theta (1 - p + pM_X(\theta))^m. \tag{7}$$

By considering the functions $f_1(y) = e^\theta (1 - p + py)^m$, $f_2(y) = \exp(\theta - \lambda + \lambda y)$, $f_3(y) = \exp(\theta - \lambda^* + \lambda^* y)$ and $f(y) = y$, and by noting that $f_1(y) \leq f_2(y)$ for all $\theta$ and $y$ and that $f_2(y) \leq f_3(y)$ for all $\theta$ and $y \geq 1$, we have

$$M_X(\theta) \leq M(\theta) \leq M^*(\theta), \tag{8}$$

where the first inequality holds for all $\theta < r_2^*$ and the second inequality holds for $0 \leq \theta < r_2^*$, and $r_2^*$ is the radius of convergence of $M^*(\theta)$, and $M(\theta)$ and $M^*(\theta)$ respectively are the smallest roots of

$$M(\theta) = \exp(\theta - \lambda + \lambda M(\theta)), \tag{9}$$
$$M^*(\theta) = \exp(\theta - \lambda^* + \lambda^* M^*(\theta)). \tag{10}$$

(Again, it can be checked that $M$ is the moment generating function for $H$.) By observing that $r_2^*$ is the value of $\theta$ at which the line $f(y) = y$ is a tangent to the curve $f(y) = \exp(\theta - \lambda^* + \lambda^* y)$, we find that $r_2^* = \lambda^* - \log \lambda^* - 1$. Further, by considering $\theta$ close to $r_2^*$, we see that $\lambda^* M^*(\theta) < 1$. Also, we shall need to bound $M''(\theta) = \frac{d^2 M}{d\theta^2}$. From (9), we have

$$M'(\theta) = M(\theta)/(1 - \lambda M(\theta))$$
$$M''(\theta) = M(\theta)/(1 - \lambda M(\theta))^3.$$

Using the fact that $\lambda M^*(\theta) \leq \lambda^* M^*(\theta) < 1$, it follows from the second inequality in (8) that for $0 \leq \theta < r_2^*$,

$$M''(\theta) \leq \frac{M^*(\theta)}{(1 - \lambda^* M^*(\theta))^3}.$$

Also, for $\theta \leq 0$, we have

$$M''(\theta) \leq \frac{1}{(1 - \lambda)^3} \leq \frac{1}{(1 - \lambda^*)^3}.$$

Thus, for any $\theta \leq (1 - \epsilon)r_2^*$ (where $\epsilon > 0$ is any fixed constant), we have

$$M''(\theta) \leq A, \tag{11}$$

where $A$ is a fixed constant (depending only on $\lambda^*$). Note that from (7) and (9), we have

$$\mathbf{E}[H_X] = \mathbf{E}[H] = 1/(1 - \lambda). \tag{12}$$

Consider next that $X_0 = 0$. For $n \geq 1$, let $\tau_n$ be the time elapsed when $X_t$ first returns to 0 for the $n$-th time. We shall obtain a concentration result for $\tau_n$ (when $n$ is large). Observe that $\tau_1$ equals $H_X$ in distribution (this is because $X_1$ has the same distribution

when $X_0 = 0$ or $X_0 = 1$) and so $\tau_n$ is distributed as a sum of $n$ independent copies of $H_X$. Hence, $\mathbf{E}[\tau_n] = n/(1 - \lambda)$. We shall use the inequalities

$$\mathbf{Pr}(\tau_n \geq A) \leq M_X(\theta)^n \exp(-A\theta),$$
$$\mathbf{Pr}(\tau_n \leq A) \leq M_X(-\theta)^n \exp(A\theta),$$

for any $\theta > 0$. As $M_X(\theta) \leq M(\theta)$ by (8), we shall bound $M(\theta)$. Using Taylor's theorem and (12),

$$M(\theta) = 1 + \theta/(1 - \lambda) + M''(\xi)\theta^2/2,$$

for some $\xi$ between 0 and $\theta$. Using (11), we have that as $\theta \to 0$,

$$M''(\xi) = O(1),$$

which implies that

$$M(\theta) = 1 + \theta/(1 - \lambda) + O(\theta^2).$$

Hence, for any $A > 0$ and small $\theta > 0$,

$$\mathbf{Pr}(\tau_n \geq n/(1 - \lambda) + An^{1/2})$$
$$\leq M(\theta)^n \exp(-n\theta/(1 - \lambda) - A\theta n^{1/2})$$
$$\leq \exp(O(n\theta^2) - A\theta n^{1/2}).$$

Also, we have for any $A > 0$ and small $\theta > 0$,

$$\mathbf{Pr}(\tau_n \leq n/(1 - \lambda) - An^{1/2})$$
$$\leq M(-\theta)^n \exp(n\theta/(1 - \lambda) - A\theta n^{1/2})$$
$$\leq \exp(O(n\theta^2) - A\theta n^{1/2}).$$

By putting $\theta = n^{-1/2}$, we have for any $A > 0$ and for large $n$

$$\mathbf{Pr}(|\tau_n - n/(1 - mp)| \geq An^{1/2}) = O(e^{-A}). \tag{13}$$

We therefore have the following lemma.

**Lemma 1** *Let $\tau_n$ be the time elapsed when $X_t$ first returns to 0 for the n-th time given that $X_0 = 0$. Then for any $A > 0$, we have as $n \to \infty$,*

$$\mathbf{Pr}(|\tau_n - n/(1 - \lambda)| \geq An^{1/2}) = O(e^{-A}).$$

**Lemma 2** *Suppose that $X_0 = n$ for any integer $n \geq 1$. Let $H_n = \min\{t \mid X_t = 0\}$. Then for any $A > 0$,*

$$\mathbf{Pr}(|H_n - n/(1 - \lambda)| \geq An^{1/2}) = O(e^{-A}). \tag{14}$$

*Also, we have for any $A > 0$ that*

$$\mathbf{Pr}(\exists t \leq H_n \text{ s.t. } X_t \geq n/(1 - \lambda) + An^{1/2}) = O(e^{-A}). \tag{15}$$

9

**Proof**  Simply observe that $H_n$ is distributed as a sum of $n$ independent copies of $H_X$, and so $H_n$ equals $\tau_n$ in distribution, which gives (14). Equation (15) follows from (14) and the fact that $X_t$ decreases by at most 1 in each transition. □

**Lemma 3** *Let $N_T$ be the number of times that $X_t$ equals 0 in the time interval $[0, T]$, given that $X_0 = O(\log^{10} T)$. Then for any $A > 0$, we have for any constant $A' > 0$ that*

$$\mathbf{Pr}(|N_T - T(1 - \lambda)| \geq AT^{1/2}) = O(e^{-A} + T^{-A'}). \tag{16}$$

**Proof**  Use $H$ to denote the minimum value of $t$ such that $X_t = 0$. Using (14) with $n = O(\log^{10} T)$, we have for any constant $A' > 0$ that

$$\mathbf{Pr}(H \geq \log^{11} T) = O(e^{-\log^6 T}) = O(T^{-A'}).$$

Hence, with error probability $O(T^{-A'})$ for any constant $A' > 0$, we may assume $N'_T \geq N_T \geq N'_{T-\log^{11} T}$ in distribution, where $N'_T$ is the number of times that $X_t = 0$ in the interval $[0, T]$ given that $X_0 = 0$. Now Lemma 1 implies that as $t \to \infty$,

$$\mathbf{Pr}(|N'_t - t(1 - \lambda)| \geq At^{1/2}) = O(e^{-A/(1-\lambda)}) = O(e^{-A}).$$

The lemma now follows by taking $t = T$ and $t = T - \log^{11} T$. □

**Lemma 4** *Suppose that $X_0 = 0$. With $\tau_1$ as defined in Lemma 1, we have for any $A > 0$, there is a constant $\rho \in (0, 1)$ and a constant $C > 0$ so that*

$$\mathbf{Pr}(\tau_1 \geq A) \leq C\rho^{-A}. \tag{17}$$

*For each $t$, let $R_t = \min\{k \geq 1 \mid X_{t+k} = 0\}$. That is, $R_t$ is the waiting time after time $t$ until the next return to 0. Then for any $A > 0$, there is a constant $\rho \in (0, 1)$ such that as $T \to \infty$,*

$$\mathbf{Pr}(\exists t \in [0, T] \ s.t. \ R_t \geq A) = O(T\rho^{-A}), \tag{18}$$

*and*

$$\mathbf{Pr}(\exists t \in [0, T] \ s.t. \ X_t \geq A) = O(T\rho^{-A}), \tag{19}$$

**Proof**  Since $\tau_1$ equals $H_X$ in distribution, we have

$$\mathbf{Pr}(\tau_1 \geq A) \leq M_X(\theta) \exp(-A\theta).$$

Inequality (17) follows by putting $\theta = r_2^*/2$. To show (18), let $S_i$ be the time elapsed between the $(i - 1)$-th and the $i$-th return to 0. That is, each $S_i$ equals $\tau_1$ in distribution. Let $N$ be the number of times that $X_t = 0$ for $t \in [0, T]$. As $N \leq T$, we have from (17) that there is a constant $\rho \in (0, 1)$ such that

$$\mathbf{Pr}(\exists i \leq T \ s.t. \ S_i \geq A) = O(T\rho^{-A}),$$

from which (18) follows easily. Inequality (19) follows from (18) and the fact that $X_t$ decreases by at most 1 in each transition. □

**Lemma 5** *Suppose that $X_0 = O(\log^{10} T)$. Then for any $A > 0$ and for large $T$, there is a constant $\rho \in (0,1)$ and a constant $C > 0$ such that*

$$\mathbf{Pr}(X_T \geq A) \leq C\rho^{-A} + O(T^{-A'}),$$

*for any constant $A' > 0$.*

**Proof** Let $H$ be as defined in the proof of Lemma 3. Note that for $t \geq H$, it follows from a simple coupling argument that $X_t$ is bounded above in distribution by a steady state chain $X_t'$ with the same transition probabilities as those of $X_t$. Hence,

$$\begin{aligned}
\mathbf{Pr}(X_T \geq A) &\leq \mathbf{Pr}(X_T \geq A \mid H \leq T)\mathbf{Pr}(H \leq T) + \mathbf{Pr}(H > T) \\
&\leq \mathbf{Pr}(X_t' \geq A)\mathbf{Pr}(H \leq T) + \mathbf{Pr}(H > T).
\end{aligned}$$

Now from (14), we have

$$\mathbf{Pr}(H > T) = O(T^{-A'}),$$

for any constant $A' > 0$. (Note that although $\mathbf{Pr}(H > T)$ should be exponentially small, our bound here will suffice for future applications.) To bound $\mathbf{Pr}(X_t' \geq A)$, we note that according to (4) and the comments that followed, the moment generating function $M_\pi(\theta)$ of $X_t'$ is properly defined for $\theta < \log r_1^*$. Hence, similar to proof of (17), there are constants $\rho \in (0,1)$ and $C > 0$ such that

$$\mathbf{Pr}(X_t' \geq A) \leq C\rho^{-A}.$$

The lemma now follows. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

Suppose that $X_0 = 0$. With $\tau_n$ defined as before, let $W_n = \sum_{t=0}^{\tau_n} X_t$. We would like to show that $W_n$ is concentrated (as $n \to \infty$). Note that

$$W_n = U_1 + \ldots + U_n,$$

where $U_1, \ldots, U_n$ are independent copies of $W_1$. To find $\mathbf{E}[W_1]$, let $W_i'$ be the sum $X_1 + \ldots + X_\tau$ given that $X_0 = i$ and $\tau$ is the first time that $X$ gets to 0. Then it can be checked that

$$\begin{aligned}
W_i' &= W_1' + (i-1)H_X + W_{i-1}', \\
W_1' &= B + W_B', \\
W_1 &= B + W_B',
\end{aligned}$$

where all equations hold in distribution only and $B$ is a binomial variable with parameters $m$ and $p$. It follows that

$$\mathbf{E}[W_i'] = i\mathbf{E}[W_1'] + \frac{i(i-1)}{2(1-mp)},$$

and so

$$\mathbf{E}[W_1] = \mathbf{E}[W_1'] = \frac{mp(2-p-mp)}{2(1-mp)^2}.$$

11

We therefore have

$$\mathbf{E}[W_n] = n\mathbf{E}[W_1] = \frac{nmp(2 - p - mp)}{2(1 - mp)^2}. \tag{20}$$

We first bound $W_1$. Using (17), we have

$$\mathbf{Pr}(\tau_1 \geq \log^2 n) = O(n^{-A}),$$

for any constant $A > 0$. (For the rest of this section, we use $A$ to denote a constant independent of $n, m, p$). Also we make no attempt to minimize the powers of $\log n$ or $\log T$.) Also, using Lemma 1, we have that

$$\mathbf{Pr}(\tau_n \geq 2n/(1 - \lambda)) = O(n^{-A}).$$

Thus, using (19), we have

$$\mathbf{Pr}(\exists t \in [0, \tau_n] \ s.t. \ X_t \geq \log^2 n) = O(n^{-A}).$$

Using the bounds on $\tau_1$ and $X_t$ and by noting that $W_1 = \sum_{t=0}^{\tau_1} X_t$, we have

$$\mathbf{Pr}(W_1 \geq \log^4 n) = O(n^{-A}), \tag{21}$$

Next let $\hat{U}_i = \min\{U_i, \log^4 n\}$, $i = 1, 2, \ldots, n$. Then we have

$$\mathbf{Pr}(\exists i \ s.t. \ \hat{U}_i \neq U_i) = O(n^{-A}). \tag{22}$$

Note that

$$\mathbf{E}[\hat{U}_i] = \mathbf{E}[W_1](1 - O(n^{-A})) + O(n^{-A} \log^4 n). \tag{23}$$

Since $\hat{U}_1, \ldots, \hat{U}_n$ are independent and bounded, we apply Hoeffding's theorem to obtain

$$\mathbf{Pr}\left(\left|\sum_{i=1}^{n} \hat{U}_i - n\mathbf{E}[\hat{U}_1]\right| \geq \sqrt{n}\log^5 n\right) \leq 2\exp(-2\log^2 n) = O(n^{-A}).$$

It therefore follows from (20), (22) and (23) that

$$\mathbf{Pr}\left(\left|W_n - \frac{nmp(2 - p - mp)}{2(1 - mp)^2}\right| \geq \sqrt{n}\log^5 n\right) = O(n^{-A}). \tag{24}$$

We now use (24) to prove the following lemma.

**Lemma 6** *Suppose that $X_0 = O(\log^{10} T)$. Let $S_T = \sum_{t=0}^{T} X_t$. Then we have*

$$\mathbf{Pr}\left(\left|S_T - \frac{Tmp(2 - p - mp)}{2(1 - mp)}\right| \geq T^{1/2}\log^6 T\right) = O(T^{-A}), \tag{25}$$

*for any constant $A > 0$.*

12

**Proof**     Let $N$ be the number of times that $X_t = 0$ for $t \in [0, T]$ and let $H$ be the minimum value of $t$ such that $X_t = 0$. Then we have

$$W_{N-1} \leq S_T \leq W_N + \sum_{t=0}^{H} X_t,$$

in distribution. Lemma 3 gives that

$$\mathbf{Pr}(|N - T(1 - mp)| \geq T^{1/2} \log^2 T) = O(T^{-A}).$$

Using (24), we have

$$\mathbf{Pr}\left(\left|W_N - \frac{Tmp(2 - p - mp)}{2(1 - mp)}\right| \geq T^{1/2} \log^5 T\right) = O(T^{-A}).$$

From Lemma 4 (see also proof of (21)), we see that $W_{N-1}$ differs from $W_N$ by at most $\log^4 T$ with probability $1 - O(T^{-A})$. Hence, we have a similar estimate for $W_{N-1}$. Using Lemma 2, we see that with probability $1 - O(T^{-A})$, we have

$$X_t = O(\log^{10} T), \quad t \leq H,$$

and

$$H = O(\log^{10}).$$

Thus, with probability $1 - O(T^{-A})$, we have

$$\sum_{t=0}^{H} X_t \leq \log^{21} T.$$

The lemma now follows.                                                                □

# 3   Proof of Theorems 1(a) and 1(c)

We shall first assume Theorem 1(b) and prove Theorem 1(c) by a monotonicity argument to show that when $c > c_3$, the probability that GUC succeeds is $o(1)$. We first consider the monotonicity argument. Suppose that we have two random instances of $k$-SAT on $n$ variables with $m$ and $m'$ clauses of size $k$ respectively. Assume $m \leq m'$. Let $N(\nu) = (N_0(\nu), N_1(\nu), N_2(\nu), N_3(\nu))$ and $\hat{N}(\nu) = (\hat{N}_0(\nu), \hat{N}_1(\nu), \hat{N}_2(\nu), \hat{N}_3(\nu))$ denote their respective states in GUC when there are $\nu$ variables whose truth values remain undetermined. We aim to show that $N(\nu) \leq \hat{N}(\nu)$ in distribution by a coupling argument. Note that the transition probabilities of $N$ are given at the end of Section 1 and that the transition probabilities of $\hat{N}$ are defined similarly with $\Delta$ replaced with $\hat{\Delta}$ and $N$ with $\hat{N}$. Note also that $N(n) = (0, \ldots, 0, m)$ and $\hat{N}(n) = (0, \ldots, 0, m')$ and so $N(n) \leq \hat{N}(n)$. We shall show inductively that if $N(\nu) \leq \hat{N}(\nu)$, then $N(\nu - 1) \leq \hat{N}(\nu - 1)$ by coupling arguments.

13

**Lemma 7** *If $N(\nu) \leq \hat{N}(\nu)$, then the chains $N$ and $\hat{N}$ can be coupled so that $N(\nu - 1) \leq \hat{N}(\nu - 1)$.*

**Proof**     Let $i \geq 1$ be the minimum integer such that $\hat{N}_i(\nu) \neq 0$. Now for $j \neq i$, $\chi_j(\hat{N}(\nu)) = 0$ and $\chi_i(\hat{N}(\nu)) = 1$. Thus, for $j \neq i$,

$$N_j(\nu) - \chi_j(N(\nu)) \leq \hat{N}_j(\nu) - \chi_j(\hat{N}(\nu)).$$

For $j = i$, we have $\hat{N}_i(\nu) \geq 1$ and note that if $N_i(\nu) = 0$ then

$$N_i(\nu) - \chi_i(N(\nu)) = 0 \leq \hat{N}_i(\nu) - 1 = \hat{N}_i(\nu) - \chi_i(\hat{N}(\nu)),$$

and that if $N_i(\nu) \geq 1$ then $\chi_i(N(\nu)) = 1$, from which we have

$$N_i(\nu) - \chi_i(N(\nu)) = N_i(\nu) - 1 \leq \hat{N}_i(\nu) - 1 = \hat{N}_i(\nu) - \chi_i(\hat{N}(\nu)).$$

Therefore, we have for all $i = 1, \ldots, k$,

$$N_i(\nu) - \chi_i(N(\nu)) \leq \hat{N}_i(\nu) - \chi_i(\hat{N}(\nu)). \tag{26}$$

Observe next that for any two binomial variables $B = B(\tau, p)$ and $\hat{B} = B(\hat{\tau}, p)$ with $\tau \leq \hat{\tau}$, we can couple $B$ and $\hat{B}$ so that

$$\begin{aligned} B &\leq \hat{B}, \\ \tau - B &\leq \hat{\tau} - \hat{B}. \end{aligned}$$

It follows from (26) that for $i = 1, \ldots, k$,

$$\Delta_{i,0}(\nu) \leq \hat{\Delta}_{i,0}(\nu), \tag{27}$$

$$N_i(\nu) - \chi_i(N(\nu)) - \Delta_{i,0}(\nu) \leq \hat{N}_i(\nu) - \chi_i(\hat{N}(\nu)) - \hat{\Delta}_{i,0}(\nu). \tag{28}$$

It follows from (27) that for $i = 1, \ldots, k$,

$$\Delta_{i,1}(\nu) \leq \hat{\Delta}_{i,1}(\nu). \tag{29}$$

Combining (28) and (29) gives that $N(\nu - 1) \leq \hat{N}(\nu - 1)$. □

**Proof of Theorem 1(c).**   For $c > c_3$, we have $c > c_3 - \epsilon$ for any $\epsilon > 0$. Now for a random instance $\mathcal{I}_\epsilon$ of 3-SAT with $\lfloor (c_3 - \epsilon)n \rfloor$ clauses and $n$ variables, Theorem 1(b) gives that the limit (as $n \to \infty$) of the probability that GUC succeeds when applied to $\mathcal{I}_\epsilon$ is arbitrarily close to 0 for sufficiently small $\epsilon > 0$. Theorem 1(c) thus follows from monotonicity. □

To show Theorem 1(a), we apply a result of Chvátal and Reed [4] which can be stated as follows. Suppose that $c < 2/3$ and consider applying algorithm SC to a random instance of 3-SAT with $n$ variables and $\lfloor cn \rfloor$ clauses. Then the probability that SC succeeds equals $1 - o(1)$ as $n \to \infty$. Theorem 1(a) now follows from the following lemma.

14

**Lemma 8** *Consider applying both GUC and SC to a random instance of $k$-SAT with $n$ variables and $m$ clauses. Then*

$$\Pr(\text{SC succeeds}) \leq \Pr(\text{GUC succeeds}).$$

**Proof**    Consider applying both SC and GUC to a random instance $\mathcal{I}$ of $k$-SAT with $n$ variables and $m$ clauses. Let $N(\nu) = (N_0(\nu), \ldots, N_k(\nu))$ and $N'(\nu) = (N_0'(\nu), \ldots, N_k'(\nu))$ denote the respective states of $\mathcal{I}$ in GUC and SC when there are $\nu$ variables whose truth values remain undetermined. Note that $N(n) = N'(n)$ initially and that the transition probabilities of $N(\nu)$ and $N'(\nu)$ are given at the end of Section 1. Note also that if $N(\nu) \leq N'(\nu)$ then $\Delta_{j,0}' \leq \Delta_{j,0}$ in distribution. Thus, by following the coupling arguments in proof of Lemma 7, we have that if $N(\nu) \leq N'(\nu)$ then $N(\nu - 1)$ and $N'(\nu - 1)$ can be coupled so that $N(\nu - 1) \leq N'(\nu - 1)$. This shows in particular that $N_0(\nu) \leq N_0'(\nu)$ in distribution, and so the lemma follows.                                                                                     $\square$

# 4    Proof of Theorem 1(b)

Assume $c \in (2/3, c_3)$. Recall that

$$f(x) = f_c(x) = \frac{3c}{4}(1 - x^2) + \log x, \quad x \in (0, 1),$$

and $c_3$ is the maximum value of $c$ such that $f(x) \leq 1$ for all $x \in (0, 1)$. Let $\alpha = \alpha(c)$ (for $c > 2/3$) be the root of the equation $f(x) = 0$ that is strictly less than 1. Note that $\alpha$ is uniquely defined and that $\alpha$ is positive. By investigating the behaviour of $f(\alpha(1 + \epsilon))$ for small $\epsilon > 0$, we see that $c\alpha^2 < 2/3$ and also if

$$\alpha_0 = \alpha + n^{-0.24}$$

then

$$nf(\alpha_0) = \Theta(n^{0.76}).$$

Note that both $\alpha n$ and $\alpha_0 n$ equal $\Omega(n)$. We shall show that if $c \in (2/3, c_3)$, then $N_2(\nu)$ can be approximated by $\nu f(\nu/n)$ as $\nu$ decreases from $n$ to $\alpha_0 n$. We shall also show that if $c$ and $\nu$ are within these ranges, then $N_3(\nu)$ can be approximated by $c\nu(\nu/n)^2$. (Thus, when $\nu = \lfloor \alpha_0 n \rfloor$, we see that $N_2(\nu) = \Theta(\nu^{0.76})$ and $N_3(\nu) \approx c\alpha_0^3 n$). These estimates enable us to find the limit of the probability that GUC succeeds.

In order to minimize subscripts, we write $W(\nu) = N_1(\nu)$, $Y(\nu) = N_2(\nu)$ and $Z(\nu) = N_3(\nu)$. We shall also consider a process $X(\nu)$ which runs alongside $N(\nu)$, and so we have a Markov chain $(N_0, W(\nu), X(\nu), Y(\nu), Z(\nu))$. The transition probabilities of $(N_0, W, Y, Z)$ are same as $N$, but those of $X$ need defining. For completeness, we write down the one-step transitions of $(W(\nu), X(\nu), Y(\nu), Z(\nu))$ below.

$$\Delta Z(\nu) = -\Delta_{3,0} - \chi_3((N_0, W, Y, Z))$$

$$\Delta Y(\nu) = \Delta_{3,1} - \Delta_{2,0} - \chi_2((N_0, W, Y, Z))$$
$$\Delta X(\nu) = \Delta_{2,1} - \chi_1((N_0, W, Y, Z)),$$
$$\Delta W(\nu) = \Delta_{2,1} - \Delta_{1,0} - \chi_1((N_0, W, Y, Z)),$$
$$\Delta N_0(\nu) = \Delta_{1,1}(\nu),$$

where

$$\Delta_{3,0} = \Delta_{3,0}(\nu) = B(Z - \chi_3((N_0, W, Y, Z)), 3/\nu),$$
$$\Delta_{3,1} = \Delta_{3,1}(\nu) = B(\Delta_{3,0}, 1/2),$$
$$\Delta_{2,0} = \Delta_{2,0}(\nu) = B(Y - \chi_2((N_0, W, Y, Z)), 2/\nu),$$
$$\Delta_{2,1} = \Delta_{2,1}(\nu) = B(\Delta_{2,0}, 1/2),$$
$$\Delta_{1,0} = \Delta_{1,0}(\nu) = B(W - \chi_1(((N_0, W, Y, Z)), 1/\nu),$$
$$\Delta_{1,1} = \Delta_{1,0}(\nu) = B(\Delta_{1,0}, 1/2).$$

The initial state of the process is $(N_0(n), W(n), X(n), Y(n), Z(n)) = (0, 0, 0, 0, \lfloor cn \rfloor)$. As the transitions of $X(\nu)$ ignores the effects of $-\Delta_{10}(\nu)$, we have $W(\nu) \leq X(\nu)$ always (which can be checked by considering the cases where $X(\nu) = W(\nu)$ and $X(\nu) > W(\nu)$). We shall see that $X(\nu)$ is a good approximation of $W(\nu)$.

We shall need the following bounds for sums of independent binomial variables. Let $B_1(\tau_1, p_1)$, $\ldots, B_k(\tau_k, p_k)$ be independent binomial variables. Write $\tau = \tau_1 + \ldots + \tau_k$ and $\bar{p} = \sum_i \tau_i p_i / \tau$. Then for $A$ satisfying $0 < A < \tau \bar{p} / 3$

$$\mathbf{Pr}\left(|B_1 + \ldots + B_k - \tau \bar{p}| \geq \sqrt{3 A \tau \bar{p}}\right) \leq 2 \exp(-A). \tag{30}$$

Also, for a binomial variable $B(\tau, p)$, we have for $u \geq e$,

$$\mathbf{Pr}(B \geq u \tau p) \leq (e/u)^{u \tau p}. \tag{31}$$

All our subsequent error probabilities regarding sums like $\sum \Delta_{3,0}$ are derived from one of the above inequalities. We shall be bounding such sums by sums of independent binomial variables. Although the variables in sums like $\sum \Delta_{3,0}$ are usually not independent, it is not difficult to show the stochastic dominance by induction and by conditioning on the outcomes of the partial sums. Also, we say that an event $\mathcal{E}$ occurs with high probability (w.h.p. for short) if

$$\mathbf{Pr}(\mathcal{E}) = 1 - O(n^{-A}), \tag{32}$$

for any constant $A > 0$. Now the events $\mathcal{E}$ usually contain bounds, involving some big $O$ terms, for random variables. In this situation, it will be clear that equations like (32) hold for any $A > 0$ by choosing sufficiently large constants (which may depend on $A$) in the big $O$ terms. We first prove the following lemma which will be useful for future inductive proofs. Note that we make no attempt to minimize the powers of $\log n$.

**Lemma 9** *Suppose that $\nu \geq \alpha_0 n$. Let $h = \lfloor n^{1/2} \rfloor$, $\nu' = \nu - h$ and $I = \{\nu' + 1, \ldots, \nu\}$. Suppose that at stage $\nu$,*

$$
\begin{aligned}
Z(\nu) &= c\nu^3/n^2 + z(n), \\
Y(\nu) &= \nu f(\nu/n) + y(n), \\
W(\nu) &= w(n) \leq \log^{10} n,
\end{aligned}
$$

*where $z(n) = o(n)$ and $y(n) = o(n^{0.76})$. Then with high probability,*

$$Z(\nu') = c\nu'^3/n^2 + O(z(n) + n^{1/4}\log n), \tag{33}$$

$$Y(\nu') = \nu' f(\nu'/n) + O(y(n) + z(n)n^{-1/2} + n^{1/4}\log n), \tag{34}$$

$$W(\nu') \leq \log^2 n, \tag{35}$$

$$\sum_{j \in I} \frac{(W(j) - 1)^+}{j} = \frac{n^{1/2} f(\nu/n)^2}{2\nu(1 - f(\nu/n))} + O(y(n)n^{-1.5} + n^{-3/4}\log^6 n). \tag{36}$$

*(The constants in the big $O$ terms are independent of $\nu$.)*

When proving the above lemma, we shall obtain the following estimates which will be useful later.

**Lemma 10** *With hypotheses of Lemma 9, we have with high probability that for all $j \in I$,*

$$Z(j) = Z(\nu) + O(n^{1/2}), \tag{37}$$

$$Y(j) = Y(\nu) + O(n^{1/2}). \tag{38}$$

*Let $\tau$ be the minimum value of $k \geq 0$ such that $W(\nu - k) = 0$, and for $j \in I$, let $\tau_j$ be the minimum value of $k \geq 1$ such that $W(j - k) = 0$. Then we have with high probability that*

$$\tau = O(w(n) + \sqrt{w(n)} \log n), \tag{39}$$

$$\tau_j \leq \log^2 n, \quad \text{for } j \leq \nu - \tau. \tag{40}$$

*Also, we have with high probability that for $j \geq \nu - \tau$,*

$$W(j) = O(w(n) + \sqrt{w(n)} \log n), \tag{41}$$

*and that for $j \leq \nu - \tau$,*

$$W(j) = O(\log^2 n). \tag{42}$$

*Note that (39-42) imply that if $w(n) = O(\log^2 n)$, then we have with high probability that for all $j \in I$,*

$$W(j) = O(\log^2 n), \tag{43}$$

$$\tau_j = O(\log^2 n). \tag{44}$$

**Proof**    We shall prove Lemma 9 and point out from where the statements in Lemma 10 follow. Note first that since $\alpha_0 n = \Omega(n)$, both $\nu$ and $Z(\nu)$ equal $\Omega(n)$. Define $\Delta'Z$ as the number of times that $Y(j) = W(j) = 0$ for $j \in I$, and $\Delta'Y$ be the number of times that $Y(j) \neq 0$ but $W(j) = 0$. Similarly, let $\Delta'W$ be the number of times that $W(j) = 0$. Therefore, we have

$$Z(\nu') - Z(\nu) = -\sum_{j \in I} \Delta_{3,0}(j) - \Delta'Z, \tag{45}$$

$$Y(\nu') - Y(\nu) = \sum_{j \in I} (\Delta_{3,1}(j) - \Delta_{2,0}(j)) - \Delta'Y. \tag{46}$$

To estimate $\sum_{j \in I} \Delta_{3,0}(j)$, we note that $\Delta_{3,0}(j)$ is bounded above in distribution by a binomial variable with parameters $Z(\nu)$ and $3/\nu'$. Thus it is not difficult to obtain that $\sum_{j \in I} \Delta_{3,0}(j)$ is bounded above by a sum of independent binomial variables, each with parameters $Z(\nu)$ and $3/\nu'$. This gives an upper bound (w.h.p.) $U_Z = O(h)$ for the sum of the variables. Since $Z(\nu') \leq Z(j) \leq Z(\nu)$, we have with high probability that $Z(j) = Z(\nu) - O(h)$, which is (37). Hence, with high probability, the distribution of $\Delta_{3,0}(j)$ is bounded below by the distribution of a binomial variable with parameters $Z(\nu) - U_Z$ and $3/\nu$. Since as $n \to \infty$,

$$\frac{3Z(\nu) - O(h)}{\nu - O(h)} = \frac{3Z(\nu)}{\nu} + O(n^{-1/2}),$$

we have with high probability that,

$$\sum_{j \in I} \Delta_{3,0}(j) = \frac{3hZ(\nu)}{\nu} + O(n^{1/4} \log n). \tag{47}$$

Similarly, we have with high probability that

$$\sum_{j \in I} \Delta_{3,1}(j) = \frac{3hZ(\nu)}{2\nu} + O(n^{1/4} \log n), \tag{48}$$

which gives us an upper bound $Y(\nu) + O(h)$ for $Y(j)$ where $j \in I$. As each $\Delta_{2,0}(j)$ is distributed as a binomial variable with parameters $Y(j) + O(1) = O(n)$ and $2/j = O(1/n)$, we have with high probability that

$$\sum_{j \in I} \Delta_{2,0}(j) = O(h).$$

Since $\Delta'Y = O(h)$, we therefore have a lower bound $Y(\nu) - O(h)$ for $Y(j)$ where $j \in I$. Thus, we have $Y(j) = Y(\nu) + O(h)$ with high probability (which is (38)). Hence, with high probability, each $\Delta_{2,0}(j)$ is bounded above and below in distribution by binomial variables with parameters $Y(\nu) + O(h)$ and $2/(\nu + O(h))$. It thus follows that

$$\sum_{j \in I} \Delta_{2,0}(j) = \frac{2hY(\nu)}{\nu} + O(n^{1/4} \log n) \tag{49}$$

with high probability. To estimate $\Delta'Z$, note first that if $\nu \leq n - n^{0.76}$ (but $\nu \geq \alpha_0 n$), then from the hypotheses in the lemma, we have

$$Y(\nu) = \Omega(n^{0.76}).$$

18

Note that during the entire time interval $I$, the number of size two clauses removed is at most $\sum_{j\in I}\Delta_{2,0}(j) + h$, which equals $O(n^{1/2})$ with high probability (using (49)). Thus the quantity $Y(j)$, for $j \in I$, is never zero and so when $\nu \leq n - n^{0.76}$,

$$\Delta' Z = 0 \tag{50}$$

with high probability. For the case where $\nu \geq n - n^{0.76}$, we consider stage $k \in I$ with $k \leq \nu - n^{0.1}$ and write $h' = \nu - k$. Then similar to (48), we have with high probability that

$$\sum_{j=k+1}^{\nu} \Delta_{3,1}(j) = \frac{3h'Z(\nu)}{2\nu}(1 + o(1)) = \frac{3c\nu^2 h'}{n^2}(1 + o(1)) = \frac{3ch'}{2}(1 + o(1)). \tag{51}$$

Also, note that $Y(n) = O(n^{0.76})$ (for $\nu \geq n - n^{0.76}$) and so similar to (49), we have that for any fixed $\epsilon > 0$,

$$\sum_{j=k+1}^{\nu} \Delta_{2,0}(j) \leq \epsilon h', \tag{52}$$

with high probability. Now in order for $Y(k) = 0$, we must have

$$\sum_{i=k+1}^{n} (\Delta_{3,1}(i) - \Delta_{2,0}(i)) \leq h'.$$

which, since $c > 2/3$ and according to (51) and (52), occurs with probability $O(n^{-A})$, for any constant $A > 0$. This shows that with high probability, $Y(k) \neq 0$ for all $k \geq \nu - n^{0.1}$. Thus with high probability, there are at most $n^{0.1}$ times when $Y(j) = 0$ (where $j \in I$). Combining this with (50), we have with high probability that

$$\Delta' Z = O(n^{0.1}). \tag{53}$$

Using (45) and (47) and (53), we have with high probability that

$$
\begin{aligned}
Z(\nu') &= Z(\nu) - \frac{3hZ(\nu)}{\nu} + O(n^{1/4}\log n) \\
&= Z(\nu)(1 - 3h/\nu) + O(n^{1/4}\log n) \\
&= Z(\nu)(\nu'/\nu)^3(1 + O(1/n)) + O(n^{1/4}\log n) \\
&= c\nu'^3/n^2 + O(z(n) + n^{1/4}\log n).
\end{aligned}
$$

This proves (33). Next, we like to estimate $\Delta' Y$ and $\Delta' W$. In view of (53), we have $\Delta' Y = \Delta' W - O(n^{0.1})$ with high probability. To estimate $\Delta' W$, we consider a process $\{X(j) \mid j \leq n\}$ with transition probabilities as defined in the beginning of this section. We also let $X(\nu) = W(\nu)$. Then as observed before, we have $W(j) \leq X(j)$ for all $j \leq \nu$. Let $\Delta' X$ be the number of times that $X(j) = 0$ for $j \in I$, and so $\Delta' W \geq \Delta' X$ (as $W(j) \leq X(j)$). Next, observe that similar to our proof of (49), we have with high probability that $\Delta_{2,1}(j)$ (for all $j \in I$) is bounded above and below in distribution by binomial variables with parameters $Y(\nu) + O(h)$ and $1/(\nu + O(h))$. Now according to the hypotheses of the lemma,

$$0 < \frac{Y(\nu) + O(h)}{\nu + O(h)} = f(\nu/n)(1 + o(1)),$$

19

which is bounded above by a constant less than 1 (since $c < c_3$). Hence, with high probability, we have that $X(j)$ (for all $j \in I$) is bounded above and below in distribution by the states of two Markov chains similar to the Markov chain described in the previous section. It therefore follows from Lemma 3 (by taking $\lambda$ there as $(Y(\nu)+O(h))/(\nu+O(h))$, $T$ there as $h$, $A$ there as $O(\log h)$) that with high probability,

$$
\begin{aligned}
\Delta' X &= h\left(1 - \frac{Y(\nu) + O(h)}{\nu + O(h)}\right) + O(h^{1/2}\log h) \\
&= n^{1/2} - \frac{Y(\nu)n^{1/2}}{\nu} + O(n^{1/4}\log n).
\end{aligned}
\tag{54}
$$

We shall show next that $\Delta'W$ and $\Delta'X$ do not differ by much. We do this by finding an estimate for

$$
\sum_{j \in I}(X(j) - W(j)),
$$

which will also be useful later. Let $\tau' = \min\{k \le \nu \mid X(j) = 0\}$ and use $\tau'_j$ to denote the minimum value of $k \ge 1$ such that $X(j - k) = 0$. Note that when $X(j) = 0$, $W(j)$ is necessarily equal to 0 (as $W \le X$). Hence whenever $\Delta_{1,0}(j) \ge 1$, its cumulative effect on $\sum W$ stops when $X$ next gets to 0. Thus,

$$
\sum_{j \in I}(X(j) - W(j)) \le \sum_{j=\tau'+1}^{\nu} X(j) + \sum_{j=\nu'+1}^{\tau'} \Delta_{1,0}(j)\tau'_j.
$$

Recall that as argued above, $X(j)$ behaves like the Markov chain $X_j$ discussed in the previous section. To estimate $\tau'$, note that if $w(n) = 0$, then $\tau' = \nu$; otherwise we apply (14) (with $n$ there as $w(n)$, and $A$ there as $O(\log n)$) to obtain that

$$
\nu - \tau' = O(w(n) + \sqrt{w(n)}\log n),
$$

holds with high probability. (Since $W(j) \le X(j)$, this gives (39).) Similarly, using (15), we have with high probability that for all $j$ between $\nu$ and $\tau$,

$$
X(j) = O(w(n) + \sqrt{w(n)}\log n),
\tag{55}
$$

from which (41) follows. Thus, with $w(n) \le \log^{10} n$, we have $\nu - \tau' = O(\log^{10} n)$ and $X(j) = O(\log^{10} n)$, from which we obtain that

$$
\sum_{j=\tau'+1}^{\nu} X(j) = O(\log^{20} n)
$$

holds with high probability. Next, for $j$ between $\nu' + 1$ and $\tau$, we have from (19) that with high probability,

$$
W(j) \le\le X(j) \le \log^2 n,
\tag{56}
$$

and so 42) follows. Next we use (18) to obtain that with high probability,

$$
\tau'_j \le \log^2 n,
\tag{57}
$$

20

from which (40) follows. (Note that strictly speaking, we have only showed that $X(j)$ can be approximated by a Markov chain $X_j$ defined in the previous section for $j \in I$. This creates a problem when estimating $\tau'_j$ for $j$ "close" to $\nu'$. However, as it can be seen easily that our previous approximations for $Z(j)$ and $Y(j)$ work for $j$ between $\nu' - \log^3 n$ and $\nu'$ also. This means that $X(j)$ can be approximated by $X_j$ for all $j$ between $\nu' - \log^3 n$ and $\nu$. As (18) gives that $\tau'_{\nu'} = O(\log^2 n)$, inequality (57) now follows from (18) too.)

Note that $\Delta_{1,0}(j)$ is a binomial variable with parameters $W(j) + O(1)$ and $1/j$. Thus it follows from (56) that $\sum_{j=\nu'+1}^{\tau'} \Delta_{1,0}(j)\tau'_j$ is bounded above by a binomial variable with parameters $O(h \log^2 n)$ and $O(1/n)$. Hence (31) gives that

$$\sum_{j=\nu'+1}^{\tau'} \Delta_{1,0}(j)\tau'_j = O(\log n)$$

with high probability. It thus follows from (56) that

$$\sum_{j=\nu'+1}^{\tau'} \Delta_{1,0}(j)\tau'_j \leq \log^3 n$$

with high probability. We thus conclude that with high probability,

$$\sum_{j \in I}(X(j) - W(j)) = O(\log^{20} n). \tag{58}$$

It follows that with high probability, we have

$$\Delta' W - \Delta' X = O(\log^{20} n).$$

This together with (54) give that with high probability,

$$\Delta' W = n^{1/2} - \frac{Y(\nu)n^{1/2}}{\nu} + O(n^{1/4} \log n). \tag{59}$$

Hence, combining (46), (48), (49), (59) and the fact that $\Delta' Y = \Delta' W - O(n^{0.1})$, we have with high probability that

$$Y(\nu') - Y(\nu) = \frac{3hZ(\nu)}{2\nu} - \frac{hY(\nu)}{\nu} - n^{1/2} + O(n^{1/4} \log n). \tag{60}$$

It follows from the hypotheses of the lemma that

$$
\begin{aligned}
& Y(\nu') \\
= \; & Y(\nu)\left(1 - \frac{n^{1/2}}{\nu}\right) + \frac{3cn^{1/2}}{2}\left(\frac{\nu}{n}\right)^2 - n^{1/2} + O(n^{1/4} \log n + z(n)n^{-1/2}) \\
= \; & f(\nu/n)(\nu - n^{1/2}) + \frac{3cn^{1/2}}{2}\left(\frac{\nu}{n}\right)^2 - n^{1/2} + O(y(n) + n^{1/4} \log n + z(n)n^{-1/2}).
\end{aligned}
$$

21

On the other hand,

$$\nu' f(\nu'/n)$$
$$= (\nu - n^{1/2})f\left(\frac{\nu}{n}\left(1 - \frac{n^{1/2}}{\nu}\right)\right)$$
$$= (\nu - n^{1/2})\left(f(\nu/n) + \left(\frac{\nu}{n}\right)^2 \frac{3cn^{1/2}}{2\nu} - \frac{n^{1/2}}{\nu}\right) + O(1)$$
$$= (\nu - n^{1/2})f(\nu/n) + \frac{3cn^{1/2}}{2}\left(\frac{\nu}{n}\right)^2 - n^{1/2} + O(1).$$

This proves (34). For (35), we use the fact that $W(\nu') \le X(\nu')$ in distribution. Since as observed previously that $X(j)$ can be approximated by a Markov chain $X_t$ defined in the previous section, inequality (35) follows from Lemma 5. It therefore remains to show (36). Note first that

$$\sum_{j \in I} \frac{(W(j) - 1)^+}{j}$$
$$= (1 + O(n^{-1/2}))\left(\frac{1}{\nu}\right)\sum_{j \in I}(W(j) - 1)^+$$
$$= (1 + O(n^{-1/2}))\left(\frac{1}{\nu}\right)\left(\sum_{j \in I}W(j) - h + \Delta'W\right)$$
$$= (1 + O(n^{-1/2}))\left(\frac{1}{\nu}\right)\left(\sum_{j \in I}W(j) + \Delta'W - n^{1/2}\right). \tag{61}$$

From (59), we have with high probability that

$$\Delta'W - n^{1/2} = -\frac{Y(\nu)n^{1/2}}{\nu} + O(n^{1/4}\log n). \tag{62}$$

Using the fact (see also discussions before (54)) that with high probability, $X(j)$ is bounded above and below by two Markov chains whose one-step transitions are governed by binomial variables with parameters $Y(\nu) + O(h)$ and $1/(\nu + O(h))$, we apply Lemma 6 with $mp$ there equals

$$\frac{Y(\nu) + O(h)}{\nu + O(h)} = \frac{Y(\nu)}{\nu} + O(n^{-1/2}),$$

and $p = 1/(\nu + O(h)) = 1/n(1 + o(1))$ and $T = h$. Thus, we have with high probability that

$$\sum_{j \in I}X(j) = \frac{n^{1/2}Y(\nu)(2 - Y(\nu)/\nu)}{2\nu(1 - Y(\nu)/\nu)} + O(n^{1/4}\log^6 n).$$

From (58), we thus have with high probability that

$$\sum_{j \in I}W(j) = \frac{n^{1/2}Y(\nu)(2 - Y(\nu)/\nu)}{2\nu(1 - Y(\nu)/\nu)} + O(n^{1/4}\log^6 n).$$

It therefore follows from (61), (62) and the above that with high probability,

$$\sum_{j \in I} \frac{(W(j) - 1)^+}{j}$$

$$= (1 + O(n^{-1/2})) \frac{n^{1/2} Y(\nu)^2}{2\nu^3 (1 - Y(\nu)/\nu)} + O(n^{-3/4} \log^6 n)$$

$$= \frac{n^{1/2} f(\nu/n)^2}{2\nu(1 - f(\nu/n))} + O(y(n) n^{-1.5}) + O(n^{-3/4} \log^6 n),$$

where the last equality follows from the fact that $Y(\nu)/\nu = f(\nu/n) + y(n)/\nu$. This proves (36). $\qquad\square$

We now make use of Lemma 9 to show Theorem 1(b). Let $h = \lfloor n^{1/2} \rfloor$ as before, and write $n_i = n - ih$ and $I_i = \{n_i + 1, \ldots, n_{i-1}\}$. Define $J$ as the greatest integer such that $n - Jh \geq \alpha_0 n$. Note first that by using induction and by applying Lemma 9 repeatedly, we have with high probability that for all $i \leq J$,

$$Z(n_i) = cn_i^3/n^2 + O(in^{1/4} \log n), \tag{63}$$

$$Y(n_i) = n_i f(n_i/n) + O(in^{1/4} \log n), \tag{64}$$

$$W(n_i) \leq \log^2 n, \tag{65}$$

where the constants in the big $O$ terms are independent of $i$. Note that since $i \leq J = O(n^{1/2})$, the error terms in the (63) and (64) are both equal to $O(n^{3/4} \log n) = o(n^{0.76})$. This implies that the values of $Z(n_i), Y(n_i)$ and $W(n_i)$ $(i \leq J)$ satisfy the hypotheses of Lemma 9, and so induction works by applying Lemma 9 repeatedly. We shall now prove the following two lemmas from which Theorem 1(b) follows immediately.

**Lemma 11**

$$\lim_{n \to \infty} \mathbf{Pr}(\text{GUC does not fail before stage } n_J) = \exp\left(-\int_\alpha^1 \frac{f(x)^2}{4x(1 - f(x))} dx\right). \tag{66}$$

**Lemma 12** *Suppose that at stage $n_J$,*

$$Z(n_J) = cn_J^3/n^2 + o(n),$$
$$Y(n_J) = n_J f(n_J/n) + o(n^{0.76}),$$
$$W(n_J) \leq \log^{10} n.$$

*Then*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{GUC creates an empty clause at and after stage } n_J) = 0. \tag{67}$$

**Proof of Lemma 11**    Note first that

$$\sum_{\nu=n_J+1}^n \frac{(W(\nu) - 1)^+}{2\nu} = \sum_{i=1}^J \sum_{j \in I_i} \frac{(W(j) - 1)^+}{2j}.$$

23

Now since from (63), (64) and (65), we have with high probability that for all $i \leq J$, the quantities $Z(n_i)$, $Y(n_i)$ and $W(n_i)$ satisfy the hypotheses of Lemma 9, we have from (36) that with high probability,

$$\sum_{j \in I_i} \frac{(W(j) - 1)^+}{2j}$$

$$= \frac{n^{1/2} f(n_{i-1}/n)^2}{4n_{i-1}(1 - f(n_{i-1}/n))} + O((i-1)n^{-1.25} \log n + n^{-3/4} \log^6 n),$$

for all $i \leq J$. Hence, we have with high probability that

$$\sum_{\nu = n_J + 1}^{n} \frac{(W(\nu) - 1)^+}{2\nu} = \sum_{i=1}^{J} \frac{n^{1/2} f(n_{i-1}/n)^2}{4n_{i-1}(1 - f(n_{i-1}/n))} + O(n^{-1/4} \log^6 n). \qquad (68)$$

Since $Jn^{-1/2} \to 1 - \alpha$ as $n \to \infty$, the sum on the right hand side of the above equation converges to

$$\int_0^{1-\alpha} \frac{f(1-x)^2}{4(1-x)(1 - f(1-x))} dx = \int_\alpha^1 \frac{f(x)^2}{4x(1 - f(x))} dx.$$

It therefore follows from (68) that for any fixed $\epsilon > 0$,

$$\left| \sum_{\nu = n_J + 1}^{n} \frac{(W(\nu) - 1)^+}{\nu} - \int_\alpha^1 \frac{f(x)^2}{2x(1 - f(x))} dx \right| \leq \epsilon, \qquad (69)$$

holds with probability tending to 1 as $n \to \infty$. Finally, as explained in (2), we have that

$$\mathbf{Pr}(\text{GUC does not fail before stage } n_J)$$

$$= \mathbf{E} \left[ \prod_{\nu = n_J + 1}^{n} \left( 1 - \frac{1}{2\nu} \right)^{(W(\nu) - 1)^+} \right]$$

$$= \mathbf{E} \left[ \exp \left( - \sum_{\nu = n_J + 1}^{n} \frac{(W(\nu) - 1)^+}{2\nu} + O \left( \sum_{\nu = n_J + 1}^{n} \frac{(W(\nu) - 1)^+}{\nu^2} \right) \right) \right].$$

The lemma therefore follows from (69) and the fact that $n_J = \Omega(n)$. $\qquad \square$

**Proof of Lemma 12**    It is useful to note that as remarked when we defined $\alpha$, the quantity $c\alpha^2$ is bounded above by a constant less than $2/3$. Note also that from the hypotheses of the lemma, we have $Z(n_J) = c\alpha^3 n(1 + o(1))$ and $Y(n_J) = o(n^{0.8})$. We consider a further $h' = \lfloor n^{0.8} \rfloor$ stages after stage $n_J$. We claim that by that stage, GUC will have arrived at a stage $n^*$ where $Y(n^*) = W(n^*) = 0$. To see this, it is not difficult to check that in these further $h'$ stages, with high probability,

(I) at most $3c\alpha^2 n^{0.8}/2(1 + o(1))$ new clauses of size 2 are created by GUC,

(II) at least $h'$ clauses of minimal sizes are removed by GUC.

(Note that (I) is similar to (48) and can therefore be proved similarly.) Since $c\alpha^2 < 2/3$ and $Y(n_J) + W(n_J) = o(n^{0.8})$, it is not possible (with high probability) to have (I) and (II) unless some of the clauses of minimal size removed are of size 3. This shows that with high

24

probability, there is $n^* \geq n_J - h'$ such that $Y(n^*) = W(n^*) = 0$. Note also that similar to (47), we have with high probability that between stages $n_J - h'$ and $n_J$, only $O(n^{0.8})$ clauses of size 3 are removed. Thus at stage $n^*$, we have with high probability that there are $Z(n^*) = c\alpha^3 n(1 + o(1))$ clauses of size three remaining, and that there are $n^* = \alpha n(1 + o(1))$ variables whose truth values remain unassigned. Since the ratio of number of size three clauses to number of variables at stage $n^*$ is strictly less than 2/3, we know from part (a) of Theorem 1 that the probability that GUC creates an empty clause at and after stage $n^*$ is $o(1)$. It therefore remains to argue that for $n$ between $n' = n_J - h'$ and $n_J$, GUC creates no empty clauses with probability tending to 1 as $n \to \infty$. To do this, note that as in (I) above, we have with high probability that

$$Y(\nu) = O(n^{0.8}),$$

for all $\nu$ between $n'$ and $n_J$. Since both $n'$ and $n_J$ equal $\Omega(n)$, we have with high probability that $Y(\nu)/\nu = o(1)$ for all $\nu \in [n', n_J]$. As indicated when showing (54), we have with high probability that for $\nu \in [n', n_J]$, $W(\nu)$ can be bounded above in distribution by a Markov chain $X_n$ defined in the previous section with one-step transitions governed by a binomial variable with parameters $O(n^{0.8})$ and $1/n'$. Using (15) and (19) and by following arguments used in showing (41) and (42), we have with high probability that for all $\nu \in [n', n_J]$, $W(\nu) \leq \log^{11} n$. This in turn gives that

$$\sum_{\nu = n'}^{n_J} \frac{W(\nu)}{\nu} = O(n^{-0.2} \log^{11} n).$$

with high probability. Since the expected number of empty clauses created at stage $\nu$ equals $O(\mathbf{E}[W(\nu)/\nu])$ (see definition of $\Delta_{1,1}$), and since there are $O(n)$ clauses, the above equation gives that the expected number of empty clauses created at stages $\nu \in [n', n_J]$ equals $o(1)$. Hence, as $n \to \infty$,

$$\mathbf{Pr}(\text{GUC creates an empty clause at stage } \nu \in [n', n_J]) = o(1). \tag{70}$$

This completes our proof of Lemma 12. □

# 5   GUC with backtracking and proof of Theorem 2

Since GUC succeeds with probability $1 - o(1)$ when $c < 2/3$, we consider only the case where $2/3 \leq c < c_3$. Note first that empty clauses can only be created by GUC when $N_1(\nu) \neq 0$. As our previous analysis shows, $N_1(\nu)$ behaves like a Markov chain in steady state with a reflecting barrier at 0. Also, given $N_1(\nu)$, the probability that GUC creates an empty clause is at stage $\nu$ is $O(N_1(\nu)/\nu)$. By allowing GUC to backtrack when it makes a "mistake", we shall see that a random instance of 3-SAT almost certainly has a satisfiable truth assignment when $c < c_3$.

Consider applying GUC to a 3-SAT problem. With $n_b > n_e$, we use $[n_b, n_e]$ to denote a "run" in which $N_1(\nu)$ is non-zero. That is, a run $[n_b, n_e]$ is such that $N_1(n_b + 1) = 0$, $N_1(k) > 0$

$(n_b \geq k > n_e)$, and $N_1(n_e) = 0$. We next describe how we allow GUC to backtrack. Recall that $N(\nu)$ is obtained from $N(\nu + 1)$ by setting a literal $x_{\nu+1}$ to 1 at stage $\nu + 1$ (using $x_\nu$ to denote the literal that is set to 1 at stage $\nu$, and recall that $x_\nu$ is a literal picked randomly from a randomly chosen clause of minimal size). Also, use $\mathcal{S}(\nu)$ to denote the set of clauses at stage $\nu$. Suppose that GUC is in a run with $N_1(n' + 1) = 0$, and $N_1(k) \geq 1$ for $k = n', n' - 1, \ldots, n''$ where $n'' \leq n'$ is the present stage. GUC then sets a literal $x_{n''}$ to 1. The backtracking is performed if the setting of $x_{n''}$ to 1 gives rise to the occurrence of two size one clauses $\{y\}$ and $\{\bar{y}\}$ for some variable $y$. If this occurs, then GUC backtracks by resetting the literals $x_{n'+1}, x_{n'}, x_{n'-1}, \ldots, x_{n''}$ to 0. We have to update the set of clauses by
(a) removing all clauses that contain $\bar{x}_k$ $(k = n' + 1, n', \ldots, n'')$ from the set $\mathcal{S}(n' + 1)$ of clauses,
(b) removing all occurrences of $x_k$ $(k = n' + 1, n', \ldots, n'')$ from clauses in the set $\mathcal{S}(n' + 1)$. Hence this new set of clauses becomes $\mathcal{S}(n'' - 1)$ and the algorithm then proceeds as before by choosing a literal $x_{n''-1}$ and set it to 1 to obtain $\mathcal{S}(n'' - 2)$. Stages $n'' - 2, n'' - 3, \ldots$ are carried out similarly as before. We call this algorithm GUCB. We say that GUCB fails if
(A) an empty clause is created in the backtracking when resetting the truth values of some literals to 0, or
(B) it has to reset the truth value of a variable more than once. That is, GUCB fails if it creates an empty clause in a stage after a backtracking and before the next time when the number of size one clauses becomes zero.

We use $\hat{N}(\nu) = (\hat{N}_0(\nu), \hat{N}_1(\nu), \hat{N}_2(\nu), \hat{N}_3(\nu))$ to denote the state of GUCB at stage $\nu$ when applied to a random instance of 3-SAT. With $n'$ and $n''$ defined as above, we claim that at stage $n'' - 1$, the set $\mathcal{S}(n'' - 1)$ of clauses remains uniformly random.

<u>Claim.</u> If $V_{n''-1}$ is the set of variables whose truth values remain unassigned at stage $n'' - 1$, then for $i = 1, 2, 3$, a size $i$ clause in $\mathcal{S}(n'' - 1)$ is equally likely to be any clause in $\mathcal{C}_i(V_{n''-1})$.
**Proof**    Let $C$ be a clause of size $s$ in $\mathcal{S}(n' + 1)$. Note that $s \geq 2$. It is clear that if $C \cap \{x_i, \bar{x}_i\} = \emptyset$ for all $i = n' + 1, n', \ldots, n''$, then $C$ is equally likely to be any clause in $\mathcal{C}_s(V_{n''-1})$. On the other hand, if $C \cap \{x_i, \bar{x}_i\} \neq \emptyset$ for some $i = n' + 1, n', \ldots, n''$, then let $j$ be the greatest value of such $i$'s. If $\bar{x}_j \in C$, then no sub-clause of $C$ is in $\mathcal{S}(n'' - 1)$ by definition of $\mathcal{S}(n'' - 1)$. If $x_j \in C$, then $C_1 = C - \{x_j\}$ is equally likely to be any clause in the set of all clauses with size $|C_1|$ made up of variables whose truth values remain unassigned immediately after stage $j$. Now since $C$ contains $x_j$, $C$ is not considered by GUCB until backtracking. During the backtracking, $C$ is removed from $\mathcal{S}(n' + 1)$ if $C$ contains $\bar{x}_i$ for some $i = j - 1, j - 2, \ldots, n''$. Otherwise $C_2 = C - \{x_{n'+1}, x_{n'}, \ldots, x_{n''}\}$ is in $\mathcal{S}(n'' - 1)$, but then $C_2$ is equally likely to be any clause of size $|C_2|$ made up of variables in $V_{n''-1}$.    $\square$

Hence the behaviour of GUCB can be analysed by considering $\hat{N}(\nu)$. As before, we shall allow GUCB to continue after empty clauses are created, that is, we allow GUCB to continue even when it fails in cases (A) and (B) above. We shall show that the probability that GUCB fails is $o(1)$. This is done by showing that the effect of backtracking on $\hat{N}$ is negligible, and that with high probability, there are at most $\log^5 n$ times when GUCB backtracks. Note that we make no attempt to minimize the powers of $\log n$ in this section.

To minimize subscripts, we write $\hat{W}(\nu)$ for $\hat{N}_1(\nu)$, $\hat{Y}(\nu)$ for $\hat{N}_2(\nu)$ and $\hat{Z}(\nu)$ for $\hat{N}_3(\nu)$. Recall

that

$$f(x) = \frac{3c}{4}(1 - x^2) + \log x, \quad x \in (0, 1).$$

The constant $\alpha$ is defined to be the unique root of $f(x) = 0$ within the range $(0, 1)$, and $\alpha_0 = \alpha + n^{-0.24}$. Also, the integer $n_J$ is defined as the greatest integer such that $n - Jh \geq \alpha_0 n$, where $h = \lfloor n^{1/2} \rfloor$. We next define some new quantities. Let $b_0 = n + 1$, $l_0 = n + 1$ and $f_0 = n + 1$. For integers $1 \leq i \leq \log^5 n_0$, if GUCB backtracks for at least $i$ times before stage $n_J$, then define $b_i, l_i, f_i$ so that $b_i$ equals the stage number at which GUCB backtracks for the $i$-th time, $l_i$ equals the greatest integer $k \leq b_i$ such that $\hat{W}(k) = 0$, and $f_i$ equals the smallest integer $k \geq b_i$ such that $\hat{W}(k + 1) = 0$; if GUCB backtracks for less than $i$ times before stage $n_J$, then define $b_i = b_{i-1}$, $l_i = l_{i-1}$ and $f_i = f_{i-1}$. (That is, $[f_i, l_i]$ is essentially a "run" corresponding to GUCB in which the backtracking takes place at stage $b_i$). We shall use induction to show that with high probability, we have for all $i \leq \log^5 n_0$ that

$$\hat{Z}(b_i - 1) = cb_i^3/n^2 + O(in^{3/4} \log n), \tag{71}$$

$$\hat{Y}(b_i - 1) = b_i f(b_i/n) + O(i^2 n^{3/4} \log n), \tag{72}$$

$$\hat{W}(b_i - 1) = O(\log^4 n), \tag{73}$$

where the constants in the big $O$ terms are independent of $i$. Note that the quantities $\hat{Z}(b_i - 1), \hat{Y}(b_i - 1), \hat{W}(b_i - 1)$ respectively are the numbers of size three, size two, size one clauses immediately after the backtracking at stage $b_i$. When proving the above using induction, it is convenient to show at the same time the following estimates that for $i \leq \log^5 n$,

$$\mathbf{Pr}(\text{GUCB creates an empty clause at stage } j \in [b_i - 1, l_i + 1]) = O(\log^8 n/n), \tag{74}$$

$$\mathbf{Pr}(\text{GUCB creates an empty clause at stage } b_{i+1}) = O(\log^6 n/n). \tag{75}$$

That (71 - 73) hold for $i = 0$ is trivial. Assume therefore that they hold for $i$, and show that (71 - 73) remain valid for $i + 1$. Note that after stage $b_i$, GUCB behaves like GUC until the next backtracking. Therefore, consider applying GUC to a random instance $\mathcal{I}$ of a satisfiability problem on $b_i - 1$ variables with $\hat{Z}(b_i - 1)$ size three clauses, $\hat{Y}(b_i - 1)$ size two clauses and $\hat{W}(b_i - 1)$ size one clauses. Use $Z(j), Y(j)$ and $W(j)$ to denote the numbers of size three, size two and size one clauses at stage $j \leq b_i - 1$. Also, for $j \leq b_i - 1$, use $\tau_j$ to denote the minimum value of $k \geq 1$ such that $W(j - k) = 0$. Note that until the next backtracking at stage $b_{i+1}$, we have $\hat{Z} = Z$, $\hat{Y} = Y$ and $\hat{W} = W$.

Note that the values of $Z, Y, W$ satisfy the hypotheses of Lemma 9. Thus, we apply (39) and (41) to obtain that with high probability,

$$b_i - l_i = O(\log^4 n), \tag{76}$$

$$W(j) = O(\log^4 n), \quad \text{for all } j \in [b_i - 1, l_i]. \tag{77}$$

We therefore have with high probability that

$$\sum_{j=b_i-1}^{l_i+1} \frac{W(j)}{j} = O(\log^8 n/n).$$

27

Hence, the expected number of empty clauses created in stages $j \in [b_i - 1, l_i + 1]$ equals $O(\log^8 n/n)$ (please refer to comments before (70)). Equation (74) now follows.

Next, we apply Lemma 9 to obtain that with high probability

$$
\begin{aligned}
Z(n') &= cn'^3/n^2 + O(in^{3/4}\log n + n^{1/4}\log n), \\
Y(n') &= n'f(n'/n) + O(in^{3/4}\log n + (i+1)n^{1/4}\log n), \\
W(n') &\leq \log^2 n,
\end{aligned}
$$

where $n' = b_i - 1 - h$. These estimates satisfy the hypotheses of Lemma 9. Therefore, if $n' \geq \alpha_0 n$, we may apply Lemma 9 repeatedly. Since we need only apply Lemma 9 at most $O(n^{1/2})$ times before we go past the stage $\lfloor \alpha_0 n \rfloor$, we have by using (37), (38), (41), (43), (40) and (44) that with high probability,

$$
\begin{align}
Z(j) &= cj^3/n^2 + O((i+1)n^{3/4}\log n), \tag{78}\\
Y(j) &= jf(j/n) + O((i+1)^2 n^{3/4}\log n), \tag{79}\\
W(j) &= O(\log^2 n) \tag{80}\\
\tau_j &= O(\log^2 n), \tag{81}
\end{align}
$$

for all $j \in [l_i, n_J]$. Note that if there are at most $i$ backtracking before stage $n_J$, then (71 - 73) remain valid for $i+1$. Otherwise, we have $l_i > f_{i+1} \geq b_{i+1} \geq n_J$ by definitions of $f_{i+1}$ and $b_{i+1}$. Therefore, using the above estimates, we have with high probability that

$$
\begin{align}
Z(b_{i+1}) &= cb_{i+1}^3/n^2 + O((i+1)n^{3/4}\log n), \tag{82}\\
Y(b_{i+1}) &= b_{i+1}f(b_{i+1}/n) + O((i+1)^2 n^{3/4}\log n), \tag{83}\\
Z(f_{i+1}+1) &= cf_{i+1}^3/n^2 + O((i+1)n^{3/4}\log n), \tag{84}\\
Y(f_{i+1}+1) &= f_{i+1}f(f_{i+1}/n) + O((i+1)^2 n^{3/4}\log n_0). \tag{85}
\end{align}
$$

Note that from (81), we have with high probability that the length of every "run" equals $O(\log^2 n)$ in the entire history when GUC is applied to a random instance $\mathcal{I}$ defined above. Thus, when GUCB backtracks at stage $b_{i+1}$, we have with high probability that GUCB need only reset the truth values of $v = O(\log^2 n)$ variables. Also, we have with high probability that

$$
f_{i+1} - b_{i+1} = O(\log^2 n). \tag{86}
$$

We next show that the backtracking does not change the numbers of size three and size two clauses by much. Note first that by (82 - 86), we have

$$
\begin{align}
Z(f_{i+1}+1) &= cb_{i+1}^3/n^2 + O((i+1)n^{3/4}\log n), \tag{87}\\
Y(f_{i+1}+1) &= b_{i+1}f(b_{i+1}/n) + O((i+1)^2 n^{3/4}\log n). \tag{88}
\end{align}
$$

Recall that in the backtracking at stage $b_{i+1}$, GUCB resets the truth values of $v = O(\log^2 n)$ variables and obtain the set of clauses at stage $b_{i+1} - 1$ by updating the set $\mathcal{S}(f_{i+1}+1)$ of clauses at stage $f_{i+1} + 1$. We next observe that in the initial set of $\lfloor cn \rfloor$ (random) clauses of size three, the expected number of clauses containing a given literal equals $O(1)$. Thus, we have with high probability that for any literal $x$, the number of clauses containing $x$ equals

$O(\log^2 n)$. Hence, with high probability, the number of size three clauses in $\mathcal{S}(f_{i+1} + 1)$ containing (at least) one of the $v$ variables is $O(\log^4 n)$. This gives that with high probability,

$$\hat{Z}(b_{i+1} - 1) - \hat{Z}(f_{i+1} + 1) = O(\log^4 n). \tag{89}$$

For size two clauses, we note first that at most $\hat{Z}(b_{i+1} - 1) - \hat{Z}(f_{i+1} + 1)$ clauses of size two are added to $\mathcal{S}(f_{i+1} + 1)$. Also, similar to (89), we have with high probability that at most $O(\log^4 n)$ size two clauses are removed from $\mathcal{S}(f_{i+1} + 1)$ in the backtracking. Therefore, we have with high probability that

$$\hat{Y}(b_{i+1} - 1) - \hat{Y}(f_{i+1} + 1) = O(\log^4 n). \tag{90}$$

Similarly, it is easy to see that with high probability, at most $O(\log^4 n)$ clauses of size 1 are created from clauses of size two and size three in $\mathcal{S}(f_{i+1} + 1)$. We thus have with high probability that

$$\hat{W}(b_{i+1} - 1) = O(\log^4 n). \tag{91}$$

The induction proof of (71 - 73) is now complete by noting that (71 - 73) follow from (87 - 91) and the fact that $\hat{Z}(f_{i+1} + 1) = Z(f_{i+1} + 1), \hat{Y}(f_{i+1} + 1) = Y(f_{i+1} + 1)$.

We next would like to show (75). Let $I = \{b_{i+1}, b_{i+1} + 1, \ldots, f_{i+1}, f_{i+1} + 1\}$ and use $V_b$ to denote the set of variables whose truth values remain unassigned immediately before stage $b_{i+1} - 1$. For $j \in I$, use $x_j$ to denote the literal that was set to 1 at stage $j$. (Note that $v = f_{i+1} - b_{i+1} + 2 = O(\log^2 n)$.) Now in the backtracking at stage $b_{i+1}$, GUCB resets these $v$ literals to 0 and update the set $\mathcal{S}(f_{i+1} + 1)$ of clauses. For $j \in I$, let $\mathcal{S}_j^{(i)}$ be the set of clauses of size $i$ in the set $\mathcal{S}(j)$ of clauses at stage $j$ containing the literal $x_j$. That is,

$$\mathcal{S}_j^{(i)} = \{C \in \mathcal{S}(j) \mid x_j \in C \text{ and } |C| = k\}.$$

Note that if $C \in \mathcal{S}_j^{(1)}$, then $C$ must come from a clause $C' \in \mathcal{S}(f_{i+1} + 1)$ where $C'$ contains a literal $\bar{x}_{j'}$, for some $j' \in I$ and $j' > j$. Thus, no clause in $\cup_{j \in I} \mathcal{S}_j^{(1)}$ can become an empty clause during backtracking. Note also that if $C \in \mathcal{S}_j^{(i)}$ ($i = 2, 3$), then the entire clause $C$ is removed from $\mathcal{S}_j$ at stage $j$, and so no sub-clause of $C$ can appear in $\mathcal{S}_{j'}^{(2)} \cup \mathcal{S}_{j'}^{(3)}$ for all $j' \in I$ and $j' < j$. Thus, if $C \in \mathcal{S}_j^{(i)}$ ($i = 2, 3$), then *during backtracking*, $C - \{x_j\}$ is equally likely to be a size $i - 1$ clause chosen from the set

$$C_{i-1}(V_b \cup \{\hat{x}_{j'} \mid j' \in I \text{ and } j' < j\}),$$

where $\hat{x}$ here denotes the variable of the literal $x$. Thus, if $C \in \mathcal{S}_j^{(i)}$ ($i = 2, 3$), then the probability that $C$ becomes an empty clause after the backtracking is $O(v/b_{i+1}) = O(\log^2 n/n)$. Note that for a clause $C \in \mathcal{S}(f_{i+1} + 1)$ to become an empty clause after backtracking, the clause $C$ must be contained in $\cup_{j \in I} \cup_{i=2,3} \mathcal{S}_j^{(i)}$. As argued in (89) and (90), the size of $\cup_{j \in I} \cup_{i=2,3} \mathcal{S}_j^{(i)}$ is $O(\log^4 n)$. Hence the probability that an empty clause is created in the backtracking at stage $b_{i+1}$ equals $O(\log^6 n/n)$. This proves (75).

It now follows from (74) and (75) that

$\mathbf{Pr}$(GUCB creates an empty clause at stages $j \in [b_i, l_i + 1]$, for some $i \leq \log^5 n$)
$= O(\log^{13} n/n).$

Therefore, it remains to show that

$$\mathbf{Pr}(\text{GUCB backtracks at least } \log^5 n \text{ times before stage } n_J) = o(1), \qquad (92)$$

and that

$$\mathbf{Pr}(\text{GUCB backtracks at and after stage } n_J) = o(1). \qquad (93)$$

To show (92), suppose that $l_i$ is given and note that GUCB behaves like GUC after each $l_i$ until the next backtracking at stage $b_{i+1}$. Note that using (78) and (79), we have with high probability that for $i \le \log^5 n$,

$$\begin{aligned}
\hat{Z}(l_i) &= cl_i^3/n^2 + O(n^{3/4}\log^6 n), \\
\hat{Y}(l_i) &= l_i f(l_i/n) + O(n^{3/4}\log^{11} n).
\end{aligned}$$

Also, $\hat{W}(l_i) = 0$. Next, consider applying GUC to a random satisfiability problem $\mathcal{I}'$ with $l_i$ variables, $Z'(l_i)$, $Y'(l_i)$ and $W'(l_i)$ clauses of size three, two and one respectively, where $Z'(l_i) \ge \hat{Z}(l_i)$, $Y'(l_i) \ge \hat{Y}(l_i)$ and $W'(l_i) \ge \hat{W}(l_i)$. Then by the monotonicity argument used in showing Theorem 1(c), we have $\hat{W}(j) \le W'(j)$ for $j \ge b_{i+1}$. Thus, if $b'$ is the minimum value of $\nu \le l_i$ such that when GUC is applied to $\mathcal{I}'$, the set of clauses at stage $\nu$ contains two clauses $\{y\}, \{\bar{y}\}$ for some $y$, then it is easy to see that $b_{i+1} \le b'$ in distribution. Next, fix a constant $c'' \in (c, c_3)$ and consider applying GUC to a random instance $\mathcal{I}''$ of 3-SAT initially with $\lfloor c''n \rfloor$ clauses of size 3 and $n$ variables. Note that by definitions of $l_i$ and $c''$, we have $l_i \ge n_J \ge \alpha_0''n$ (where $\alpha_0''$ is defined as $\alpha_0$ but with $c$ replaced by $c''$). Thus, we apply Lemmas 9 and 10 to obtain that with high probability, the numbers of size three, size two and size one clauses with respect to $\mathcal{I}''$ satisfy that for $i \le \log^5 n$,

$$\begin{aligned}
Z''(l_i) &= c''l_i^3/n^2 + O(n^{3/4}\log^6 n), \\
Y''(l_i) &= l_i g(l_i/n) + O(n^{3/4}\log^{11} n), \\
W''(l_i) &= O(\log^2 n),
\end{aligned}$$

where $g(x) = 3c''(1 - x^2)/4 + \log x$. Let $N''$ be the number of stages $\nu$ before $n_J$ such that in applying GUC to $\mathcal{I}''$, the set of clauses at stage $\nu$ contains two clauses $\{y\}, \{\bar{y}\}$ for some $y$. Since $Z''(l_i) \ge \hat{Z}(l_i)$, $Y''(l_i) \ge \hat{Y}(l_i)$ and $W''(l_i) \ge \hat{W}(l_i)$ with high probability, it follows (by considering the waiting times $b'$ defined above) that

$$\begin{aligned}
&\mathbf{Pr}(\text{GUCB backtracks at least } \log^5 n \text{ times before stage } n_J) \\
&\le \mathbf{Pr}(N'' \ge \log^5 n) + o(1).
\end{aligned}$$

Using (43), we see that when GUC is applied to $\mathcal{I}''$, we have with high probability that for all $j \ge n_J$, the number $W''(j)$ of size one clauses at stage $j$ is $O(\log^2 n)$. Therefore, the probability that there is a contribution to $N''$ at stage $j$ equals $O(\mathbf{E}[W''(j)^2/j])$. Since $W''(j) = O(n)$, we now have

$$\mathbf{E}[N''] = O(\log^4 n).$$

It therefore follows that

$$\mathbf{Pr}(N'' \ge \log^5 n) = O(1/\log n).$$

This shows (92).

To show (93), we have from (78 - 80) again that with high probability

$$
\begin{aligned}
\hat{Z}(n_J) &= cn_J{}^3/n^2 + O(n^{3/4}\log^6 n), \\
\hat{Y}(n_J) &= n_J f(n_J/n) + O(n^{3/4}\log^{11} n), \\
\hat{W}(n_J) &= O(\log^4 n).
\end{aligned}
$$

(Note that in the unlikely event where $n_J \in [b_i - 1, l_i]$ for some $i$, we may apply (71 - 73) and (76 - 77) to obtain the above estimates at stage $n_J$.) These values of $\hat{Z}, \hat{Y}, \hat{W}$ satisfy the hypotheses of Lemma 12. Thus, we obtain (93) from Lemma 12. Our proof of Theorem 2 is thus complete.

# 6 Proof of Theorem 3

We shall only give a sketch proof here. Consider SC when applied to a random instance of $k$-SAT with $n$ variables and $m = \lfloor cn \rfloor$ clauses. We restrict our attention to

$$
c > \left(\frac{k-1}{k-3}\right)\frac{k-1}{k-2}\frac{2^{k-3}}{k},
$$

for otherwise SC succeeds with probability $1 - o(1)$ (see Chvátal and Reed [4]). Let $q_i(\nu)$ be the probability that a randomly selected clause from $C_k(V_n)$ is of size $i$ immediately before stage $\nu$. It is not difficult to check that for $i = 3, \ldots, k$,

$$
q_i(\nu) = \frac{\binom{n-\nu}{k-i}\binom{\nu}{i}}{\binom{n}{k}}2^{-(k-i)}.
$$

Let $N_i'(\nu)$ be the number of size $i$ clauses at stage $\nu$. The above equation implies that with high probability, we have for $i = 3, \ldots, k$ that

$$
N_i'(\nu) = \binom{k}{i}\frac{cn}{2^{k-i}}(\nu/n)^i(1 - \nu/n)^{k-i} + O(n^{1/2}\log n), \tag{94}
$$

whenever $\nu = \Omega(n)$. This gives a fairly accurate estimate for $N_3'(\nu)$ in particular.

Fix a (small) constant $\epsilon > 0$. Recall that $\beta_1$ is the largest root of the equation

$$
p_3(x) = \binom{k}{3}cx^2(1 - x)^{k-3}2^{-(k-3)} = 2/3.
$$

Let $\beta_1' = \beta_1 + \epsilon$ and $\beta_1'' = \beta_1 - \epsilon$. Note that $N_1'(\nu - 1) + N_2'(\nu - 1) - N_1'(\nu) - N_2'(\nu)$ is bounded above by

$$
\begin{cases}
\Delta_{3,1}'(\nu), & \text{if } N_1'(\nu) + N_2'(\nu) = 0, \\
\Delta_{3,1}'(\nu) - 1, & \text{otherwise,}
\end{cases}
$$

31

where $\Delta'_{3,1}(\nu)$, defined in Section 1, is the number of new size 2 clauses created at stage $\nu$. Since $\Delta'_{3,1}(\nu)$ is a binomial variable with parameters $N'_3(\nu)$ and $3/2\nu$ and since for $\nu \geq \beta'_1 n$,

$$N'_3(\nu)\frac{3}{2\nu} = \frac{3c}{2}\binom{k}{3}(\nu/n)^2(1 - \nu/n)^{k-3}2^{-(k-3)} + O(n^{-1/2}\log n) < 1$$

with high probability, it follows from Lemma 4 (see also proof of (43)) that for $\nu \geq \beta'_1 n$,

$$N'_1(\nu) + N'_2(\nu) = O(\log^2 n)$$

with high probability. This gives an upper bound for $N'_2(\nu)$ which in turn gives that with high probability,

$$\sum_{n \geq \nu \geq \beta'_1 n} N'_1(\nu) = O(\log^2 n).$$

It therefore follows as in (3) that

$$\lim_{n \to \infty} \mathbf{Pr}(\text{SC fails at or before stage } \beta'_1 n) = 0. \tag{95}$$

Furthermore, for $\nu$ between $\beta''_1 n$ and $\beta'_1 n$, it is not difficult to obtain that there is $\gamma_1(\epsilon)$ which tends to 0 as $\epsilon \to 0$ such that

$$N'_2(\nu) \leq \gamma_1(\epsilon)n \tag{96}$$

with high probability. This gives an upper bound for $N'_1(\nu)$ and it is not difficult to obtain in a similar (but simpler) fashion as our proof of Theorem 1(b) that there is $\gamma_2(\epsilon)$ where $\gamma_2(\epsilon) \to 0$ as $\epsilon \to 0$ such that

$$\lim_{n \to \infty} \mathbf{Pr}(\text{SC fails at a stage between } \beta''_1 n \text{ and } \beta'_1 n) \leq \gamma_2(\epsilon). \tag{97}$$

Suppose we allow SC to have limited backtracking (as in GUC described in the previous section). Then in view of (95) and (97), the theorem follows from the following lemma.

**Lemma 13** *For all small $\epsilon > 0$,*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{SCB fails at or after stage } \beta''_1 n) = 0.$$

We do not prove Lemma 13. Instead, we give a sketch proof of Lemma 14 below. ($c < c_k$ means that Lemmas 13 and 14 can be proved similarly.)

**Lemma 14** *Let $n_0 = \lfloor \beta_1 n \rfloor$ and $V$ be a set of $n_0$ variables. Let $\mathcal{I}$ be a random formula with $\hat{N}_i(n_0)$ clauses of size $i$, where for $i = 3, \ldots, k$,*

$$\hat{N}_i(n_0) = \binom{k}{i}cn(n_0/n)^i(1 - n_0/n)^{k-i} + O(n^{1/2}\log n)$$

*and $\hat{N}_i(n_0) = 0$ for $i = 0, 1, 2$. Each size $i$ clause in $\mathcal{I}$ is chosen at random (with equal probability) and independently from $C_i(V)$. Then*

$$\lim_{n \to \infty} \mathbf{Pr}(\text{SCB, applied to } \mathcal{I}, \text{ fails at or after stage } n_0) = 0.$$

This lemma can be proved in a way similar to our proof of Theorem 2. The key point is that when SC (without backtracking) is applied to $\mathcal{I}$, we can follow our proof of (34) to obtain an estimate for the number $N_2'(\nu)$ of clauses of size two. Indeed, if $h = \lfloor n^{1/2} \rfloor$, $n_i = n_0 - ih$, $I_i = \{n_i + 1, \ldots, n_{i-1}\}$ and $J$ is the greatest integer such that $n_0 - Jh \geq \alpha n + n^{0.76}$ where $\alpha$ is defined later, then we have with high probability that

$$N_2'(n_i) = \frac{kcn_i}{2^{k-1}} \left[ (1 + (k-2)n_i/n)(1 - n_i/n)^{k-2} - (1 + (k-2)\beta_1)(1 - \beta_1)^{k-2} \right]$$
$$+ n_i \log(n_i/(\beta_1 n)) + O(in^{1/4} \log n), \tag{98}$$

which can be proved using induction and difference equations as in Lemma 9. Intuitively, the above equation can be obtained as follows. Let

$$p_2(x) = \frac{kc}{2^{k-1}}(1 + (k-2)x)(1-x)^{k-2} + \log x.$$

Note that $p_2(x) - p_2(\beta_1)$ is an approximaton to $N_2'(\lfloor xn \rfloor)/\lfloor xn \rfloor$ according to (98). We define $\alpha < \beta_0$ as the smallest number so that $p_2(x) - p_2(\beta_1) = 0$. Note also that

$$\frac{dp_2}{dx} = \frac{1}{x}\left(-\frac{3}{2}p_3(x) + 1\right). \tag{99}$$

Thus $p_2(x)$ is maximized when $x = \beta_0$. Note that

$$p_2(\beta_0) - p_2(\beta_1) = \frac{1}{(k-1)(k-2)}\left(\frac{1}{\beta_0^2} + \frac{k-3}{\beta_0} - \frac{1}{\beta_1^2} - \frac{k-3}{\beta_1}\right) + \ln(\beta_0/\beta_1),$$

which is less than 1 according to the hypothesis of the theorem. Thus, taking (98) as induction hypothesis, we see that $N_2'(n_i)/n_i$ is, with high probability, at most a constant which is less than 1. This means that we can apply the results in Section 2 to approximate $N_1'(\nu)$, and in particular obtain that (see $\pi_0$ before (4))

$$\Pr(N_1'(\nu) = 0) \approx 1 - N_2'(\nu)/\nu.$$

This shows that

$$\mathbf{E}[N_2'(\nu - 1) - N_2'(\nu)] \approx \mathbf{E}[\Delta_{3,1}'(\nu) - \Delta_{2,0}'(\nu)] - \Pr(N_1'(\nu) = 0)$$
$$\approx \frac{3}{2\nu}\mathbf{E}[N_3'(\nu)] - \frac{1}{\nu}\mathbf{E}[N_2'(\nu)] - 1.$$

Putting $\varphi(x) = \mathbf{E}[N_2'(\lfloor xn \rfloor)/\lfloor xn \rfloor]$, we have for small $h > 0$ that

$$\varphi(x - h) - \varphi(x)$$
$$\approx (1 + h/x + O(h^2))\frac{1}{xn}\mathbf{E}[N_2'(\lfloor xn - hn \rfloor)] - \frac{1}{xn}\mathbf{E}[N_2'(\lfloor xn \rfloor)]$$
$$\approx \frac{1}{xn}\left(\mathbf{E}[N_2'(\lfloor xn - hn \rfloor)] - \mathbf{E}[N_2'(\lfloor xn \rfloor)]\right) + \frac{h}{x^2 n}\mathbf{E}[N_2'(\lfloor xn \rfloor)] + O(h^2)$$
$$\approx \frac{h}{x}\left(\frac{3}{2}p_3(x) - 1\right).$$

So $\varphi(x)$ should stay close to the solution of the differential equation (99). The induction proof of (98) is completed by showing that $N_2'(n_{i+1}) - N_2'(n_i)$ is close to its mean.

It can be shown that the Claim in Section 5 remains true for SCB when applied to $\mathcal{I}$. That is, the set of clauses after each (limited) backtracking remains uniformly random. Therefore, our proof of (92) and the statement before it can be extended to show that

$$\mathbf{Pr}(\text{SCB, applied to } \mathcal{I}, \text{ fails at a stage between } n_J \text{ and } n_0) = o(1). \tag{100}$$

It therefore remains to show that

$$\mathbf{Pr}(\text{SCB, applied to } \mathcal{I}, \text{ backtracks at and after stage } n_J) = o(1). \tag{101}$$

Proving (101) requires a result similar to Lemma 12. Since the backtracking in SCB does not change $N_i'(\nu)$ by much, we have in particular estimates for $\hat{N}_i(n_J)$ (similar to those given in (94) and (98)). Thus as in the proof of Lemma 12, there is (with high probability) $n^* \approx \alpha n$ such that $\hat{N}_1(n^*) = \hat{N}_2(n^*) = 0$ and that for $i = 3, \ldots, k$ and for $\nu \leq n^*$, $\hat{N}_i(\nu)$ can be approximated by estimates similar to those given in (94). Note that for $\nu < n^*$, $N_3(\nu)/\nu$ is less than a constant which is less than $2/3$. Thus similar to (95) and (97), we have (101).

# References

[1] A.Z. Broder, A.M. Frieze and E. Upfal, On the satisfiability and maximum satisfiability of random 3-CNF formulas, to appear in *SODA 1993.*

[2] M.T. Chao and J. Franco, Probabilistic analysis of two heuristics for the 3-satisfability problem, *SIAM Journal on Computing* 15 (1986) 1106-1118.

[3] M.T. Chao and J. Franco, Probabilistic analysis of a generalization of the unit-clause literal selection heuristics for the $k$ satisfiabiable problem, *Information Science* 51 (1990) 289-314.

[4] V. Chvátal and B. Reed, Mick gets his (the odds are on his side), *Proceedings of the 33rd IEEE Symposium on Foundations of Computer Science*, (1992) 620-627.

[5] M. Davis and H. Putnam, A computing procedure for quantification theory, *Journal of the ACM* 7 (1960) 201-215.

[6] A. Goerdt, A threshold for unsatisfiability, to appear in *17th International Symposium on Mathematical Foundations of Computer Science*, Prague, Czechoslovakia, August 1992.

[7] A. Goldberg, Average case complexity of the satisfiability problem, *Proceedings of 4th Workshop on Automated Deduction*, (1979) 1-6.

[8] T. Larabee, Evidence for the satisfiability threshold for random 3CNF formulas.

[9] D. Mitchell, B. Selman and H. Levesque, Hard and easy distributions of SAT problems.

34