



## Random 3-SAT: The Plot Thickens

CRISTIAN COARFA

*Department of Computer Science, Rice University, 6100 S. Main St MS 132, Houston TX 77005-1892*

DEMETRIOS D. DEMOPOULOS

*Department of Computer Science, Rice University, 6100 S. Main St MS 132, Houston TX 77005-1892*

ALFONSO SAN MIGUEL AGUIRRE

*Department of Computer Science, Rice University, 6100 S. Main St MS 132, Houston TX 77005-1892*

DEVIKA SUBRAMANIAN

*Department of Computer Science, Rice University, 6100 S. Main St MS 132, Houston TX 77005-1892*

MOSHE Y. VARDI

*Department of Computer Science, Rice University, 6100 S. Main St MS 132, Houston TX 77005-1892*

**Abstract.** This paper presents an experimental investigation of the following questions: how does the average-case complexity of random 3-SAT, understood as a function of the order (number of variables) for fixed density (ratio of number of clauses to order) instances, depend on the density? Is there a phase transition in which the complexity shifts from polynomial to exponential in the order? Is the transition dependent or independent of the solver? Our experiment design uses three complete SAT solvers embodying different algorithms: GRASP, CPLEX, and CUDD. We observe new phase transitions for all three solvers, where the median running time shifts from polynomial in the order to exponential. The location of the phase transition appears to be solver-dependent. GRASP shifts from polynomial to exponential complexity near the density of 3.8, CPLEX shifts near density 3, while CUDD exhibits this transition between densities of 0.1 and 0.5. This experimental result underscores the dependence between the solver and the complexity phase transition, and challenges the widely held belief that random 3-SAT exhibits a phase transition in computational complexity very close to the crossover point.

### 1. Introduction

The last decade has seen an intense focus on the complexity of randomly generated combinatorial problems. This interest was stimulated by the discovery of a fascinating connection between the *density* of combinatorial problems and their computational complexity, see [12, 44]. A problem that has received a lot of attention in this area is the *3-satisfiability problem* (3-SAT), a paradigmatic combinatorial problem that is important also for its own sake. An instance of 3-SAT consists of a conjunction of clauses, each one a disjunction of three literals. The goal is to find a truth assignment that satisfies all clauses. The density of a 3-SAT instance is the ratio of the number of clauses to the number of Boolean variables. We call the number of variables the *order* of the instance. Clearly, a low density suggests that the instance is under-constrained, and therefore is likely to be satisfiable, while a high density suggests that the instance is over-constrained and is unlikely to be satisfiable. Experimental research [17, 44] has shown that for ratio below (roughly) 4.26, the probability of satisfiability of a random 3-SAT instance goes to 1 as the order increases, while for ratio above 4.26 the probability goes to 0. At 4.26, the

probability of satisfiability is near 0.5. This satisfiability threshold density has been called the *crossover point*. Theoretically establishing the density at the crossover point is difficult, and is the subject of continuing research, cf. [1, 21, 22].

The experiments in [17, 44], which applied algorithms based on the so-called *Davis-Logemann-Loveland method* (abbr., DLL method) (a depth-first search with unit propagation [20]), also show that the density of a 3-SAT instance is intimately related to its computational complexity. Intuitively, under-constrained instances are easy to solve, as a satisfying assignment can be found fast, and over-constrained instances are also easy to solve, as all branches of the search terminate quickly. Indeed, the data displayed in [17, 44] demonstrate a peak in running time essentially at the crossover point. Using finite-size scaling techniques, [35] demonstrated a *phase transition* at the crossover point, viz., a marked qualitative change in the structural properties of the problem. This pattern of computational behavior with a peak at the crossover point is called the *easy-hard-easy* pattern in [43].

This picture, however, is quite simplistic for various reasons. First, the terms “easy” and “hard” do not carry any rigorous meaning. The computational complexity of a problem is typically measured by its *scalability*, that is, its growth as a function of the input size. Thus, one studies computational complexity on an infinite collection of instances. The easy-hard-easy pattern, however, is observed when the order is fixed while the density varies, but once the order is fixed, there are only finitely many possible instances. For that reason, theoretical analyses of the random 3-SAT problem focus on collections of fixed-density instances, rather than on collections of fixed-order instances, cf. [3]. Second, in the context of a concrete application, e.g., bounded model checking [4], planning [33], or scheduling [18], it is typically the order that tends to grow while the density stays fixed, for example, as we search for longer and longer counterexamples in bounded model checking. Thus, experimental results that vary density while fixing the order tell us little about the computational complexity of 3-SAT in such settings. Finally, it is not clear where the boundaries between the so-called “easy”, “hard”, and “easy” regions are. A widely held belief [43, 44] is that random 3-SAT problems are “hard” only for densities very close to the crossover point. Much work has therefore focused on explaining the “jump” in computational complexity around the crossover point using finite-size scaling [43] and backbones [41]. This alleged “jump”, however, has not been documented experimentally. In fact, there is almost no experimental work that studies how the running time of a SAT solver varies as a function of the order for fixed-density instances (a few results of this nature, though not a systematic study, are reported in [16, 17, 43]). Further, the experiments reported in [17, 44] are based solely on DLL algorithms. While these are indeed the most popular algorithms for the satisfiability problem, one cannot draw conclusions about the inherent and practical complexity of random 3-SAT based solely on experiments using these algorithms. We may observe different phenomena by experimenting with SAT solvers that embody different algorithms.

The goal of our experimental algorithmic research reported here is to determine how the average-case complexity of random 3-SAT, understood as a function of the order for fixed density instances, depends on the density. Is there a phase transition in which the complexity shifts from polynomial to exponential? Is such a transition dependent or independent of the solver?

To explore these questions, we set out to obtain a good coverage of an initial quadrangle of the two-dimensional  $d \times n$  quadrant, where  $d$  is the density and  $n$  is the order. We explored the range  $0 \leq d \leq 15$ . We attempted to maximize the order of the sampled instances, given our resource constraints. We used three different SAT solvers, embodying different underlying algorithms. GRASP ([vinci.inesc.pt/~jpms/grasp/](http://vinci.inesc.pt/~jpms/grasp/)) is based on the DLL method, but it augments the search with a conflict-analysis procedure that enables it to backtrack non-chronologically and record the causes of conflict. Experimental results [38] show that GRASP is very efficient for a large number of realistic SAT instances, and it has proven to be a very effective SAT solver in the context of automated hardware design [42]. The CPLEX MIP Solver is a commercial optimizer for linear-programming problems with integer variables ([www.cplex.com](http://www.cplex.com)). It employs a branch-and-bound technique using linear-programming relaxations that can be complemented with the dynamic generation of cutting planes. While branch-and-bound is related to depth-first-search, the cutting-planes technique is more powerful than resolution [30]. CUDD ([bessie.colorado.edu/~fabio/CUDD](http://bessie.colorado.edu/~fabio/CUDD)) implements functions to manipulate Reduced Ordered Binary Decision Diagrams (ROBDDs), which provide an efficient representation for Boolean functions [9]. Unlike GRASP and CPLEX, CUDD does not search for a single satisfying truth assignment. Rather, it constructs a compact symbolic representation of the set of satisfying truth assignments and then checks whether this set is nonempty. Uribe and Stickel [50] compared ROBDDs with the DLL method for SAT solving, concluding that the methods are incomparable, and that ROBDDs dominate the DLL method on many examples. Recent work by Groote and Zantema proved the incomparability of ROBDDs and resolution [28]. ROBDDs have proven in the 1990s to be very effective in the context of hardware verification [10, 31].

Our aim was not to directly compare the performance of the different solvers to determine which one has the “best” performance, but rather to understand their behavior in the  $d \times n$  quadrant in order to make qualitative observations on how the complexity of random 3-SAT is viewed from different algorithmic perspectives. It is important to note that the algorithms used in GRASP, CPLEX, and CUDD do not explicitly refer to the density of the input instances. Thus, a qualitative change in the behavior of the algorithm, as a result of changing the density, indicates a genuine structural change in the SAT instances from the perspective of the algorithm.

In analyzing our experimental results we focus on measuring the *median* running time as a function of the order for a set of instances of fixed density.<sup>1</sup> This gives us a measure of the running time of the algorithm for that density. Our findings show that for GRASP and CPLEX the easy-hard-easy pattern is better described as an—easy—hard—less-hard pattern, where, as is the standard usage in computational complexity theory, “easy” means *polynomial time* and “hard” means *exponential time*. When we start with low-density instances and then increase the density, we go from a region of polynomial running time, to a region of exponential running time, where the exponent first increases and then decreases as a function of the density. Thus, we observe at least *two* phase transitions as the density is increased: for GRASP, a transition at around density 3.8 from polynomial to exponential running time and a transition at around density 4.26 (the crossover point) from an increasing exponent to a decreasing exponent. The region between 3.8 and 4.26 is also characterized by the prevalence of very hard instances, the so called “heavy-tail

phenomenon”, cf. [25, 29, 39, 43]. Our results indicate one or more phase transitions in this region, where the ratio of the mean-to-median running time peaks. For CPLEX we also observe another phase transition at around density 1.7 from linear running time to quadratic running time. Note that we use the term “phase transition” in a somewhat liberal sense. The phase transitions that we observed involve various measures: a change in the degree of the polynomial characterizing the running time, a change in running time from polynomial in the order to exponential in the order, and a change in the direction of the heavy-tail phenomenon. All these suggest marked qualitative changes in structural properties of the random 3-SAT instances.

A very different picture emerges for CUDD. Here the algorithm is exponential (in both time and space) for densities between 0.5 and 15. There is, however, no peak around the crossover point and no heavy-tail phenomenon was observed. We observed, however, a peak in the size of the final BDDs constructed by the algorithm at around density 2, indicating a phase transition at around this density. At a very low density (0.1) we did observe polynomial (cubic) behavior, which suggests that another phase transition is “lurking” between densities 0.1 and 0.5.

There are two conclusions that can be drawn from our experiments. First, the “phase transition” in average case computational complexity of random 3-SAT should not be identified with the “phase transition” in satisfiability. The sharp shift of average case complexity from polynomial to exponential occurs well before the crossover point. This implies that explanations for shifts in computational complexity cannot center around phenomena observed at the crossover point [41, 43]. Second, unlike earlier predictions (cf. [15, 36]), phase transitions in average-case complexity (unlike the one for satisfiability) are not solver-independent. This implies that any theory attempting to explain the sharp shift in computational complexity must take the characteristics of the solver into account, as in [2].

## 2. Related Work

The fact that the “easy-hard-easy” pattern is quite simplistic is known, though rather under-emphasized, cf. [2, 45]. For example, in the high-density region (above density 5.2), an exponential lower bound on the length of resolution proofs is proved in [13]. This entails an exponential lower bound on the running time of DLL algorithms, implying that the high-density region can, at best, be described as “less hard”. (Note that this lower bound does not apply to algorithms that are based on cutting planes or ROBDDs [15, 28, 30].) It is also known that the probability crossover is not the only phase transition involving random 3-SAT and that phase transitions can be solver dependent. In [47], the authors experimentally demonstrated a change from exponentially fast to power law relaxation at around density 3. In [8, 40], the authors proved linear median running time of the pure-literal algorithm at the low-density region (below 1.63) and showed a phase transition at 1.63 for this algorithm. In [23], the authors proved a linear median running time for the GUC heuristic for low-density instances and showed a phase transition near density 3 for this algorithm.

These latter results indicate that the low-density region is indeed in some sense “easy”, but they do not establish that complete SAT solvers have polynomial median

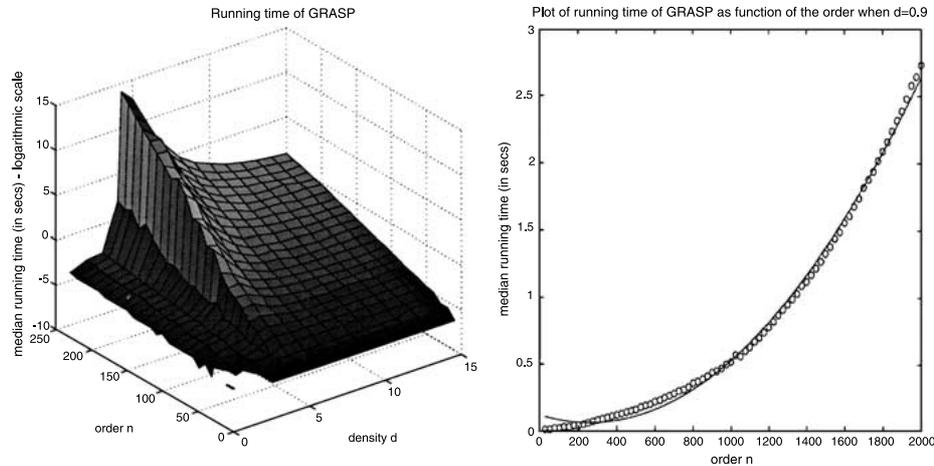


Figure 1. GRASP – (left) 3-D Plot of median running time, and (right) median running time for density 0.9 as a function of the order of the instances. A quadratic function fits these points better (with an  $r^2 > 0.98$ ) than an exponential function.

running time in this region. The analytical results of Franco and his collaborators suggest that in this region we might expect a polynomial median running time for certain heuristic algorithms, cf. [11], but they do not prove it definitively. In [16], the authors reported linear median running time of Tableau, their SAT solver, for densities 1, 2, and 3, and an exponential median running time for densities 4.26 and 10. In [43], the authors reported linear median running time of their DLL SAT solver for density 3, and an exponential median running time for density 4.26. None of these papers, however, systematically explores the dependence on the density of the running time as a function of the order.

The performance of integer-programming algorithms on random SAT instances is studied in [30], but the author did not systematically study how the running time depends on the density of the instances. Similarly, in [7, 27] the authors studied the behavior of ROBDDs on random SAT instances, but did not study how this behavior varies as a function of the density.

While we focus in this paper on the study of collections of fixed-density instances, it would also be interesting to study the behavior of SAT solvers on instances where the order and the density vary simultaneously; for example, the density may increase together with the order. For DLL solvers, the results in [3] show that unless the density increases linearly with the order we should still expect to see exponential running time. Indeed, if one considers a logarithmic increase of the density as a function of the order, then our data (e.g., using Figure 1) shows that the median running time for GRASP is still exponential.

While the finite-size scaling studies in [26, 43] do aim to explore how both the density and the order affect the running time of SAT solvers, they do not reveal the same detailed picture that emerges from our experiments on the density-order quadrant. First, the finite-size scaling studies for SAT are limited to DLL solvers. Second, finite-size scaling studies show a very good fit only around the crossover point, but the fit gets worse as the scale

value gets further from it. This makes it very difficult to draw conclusions on the dependence on the order for fixed-density instance in regions with density far below the crossover point. For example, it is not at all clear how one can obtain linear-time behavior at density 3 reported in [43] from their normalized and rescaled results. Finally, the fact that the running time of DLL is exponential in the high-density region and polynomial in the low-density region makes it rather unlikely that the scale factor observed in [43] applies anywhere but very near the crossover point. We elaborate on this point below.

Beame et al. [3] showed that in the high-density region, a certain variant of DLL terminates with high probability within time  $2^{an/d + b \log n}$ , where  $a$  and  $b$  are some positive constants,  $n$  is the order and  $d$  is the density of the instance. This matches the exponential lower bound of [13]. Thus, the *normalized* running time (i.e., the ratio of the running time to the running time at the crossover point) is  $2^{an(1/d - 1/t)}$ , where  $t$  is the crossover density. Thus, in the high-density region the finite-size scale factor is  $n(1/d - 1/t)$ . More generally, let the running time of a SAT solver in the high-density region be  $2^{an^b/d^c}$ , where  $a, b, c$ , and  $d$  are positive constants. Then the scale factor will be  $n^b(1/d^c - 1/t^c)$ . Note that in the high-density region, order and the density have opposing effects. Increasing the order increases the running time, but increasing density decreases the running time. The scale factor of  $n^\alpha(d/t - 1)$  proposed in [43] does not reflect this opposition, as it increases with both  $n$  and  $d$ , which explains why the study in [43] shows a very good fit only around the crossover point.

On the other hand, our experiments, as well as other experiments mentioned above, provide evidence for the fact that in the low-density region the running time of DLL solvers is polynomial, i.e.,  $f(d)n^e$ , where  $f$  is some function and  $e$  is a positive constant. Thus, the normalized running time is  $f(d)n^e/2^{an^b/t^c}$ . Therefore, the scale factor  $n^\alpha(d/t - 1)$  of [43] does not appear to be a reasonable scale factor in this region. For the low-density region, it makes more sense to normalize the running time with respect to some other threshold  $s$  in the low-density region. The normalized running time is then  $f(d)/f(s)$ , which implies that  $d$  is the appropriate scale factor in this region. The bottom line regarding finite-size scaling is that running times in the low- and high-density regions are very different functions of density and order. Expecting the same scale factor to work in both regions is unrealistic.

As noted above, since the sharp shift of average-case running time from polynomial to exponential is solver dependent and occurs well before the crossover point for all the solvers we tested, explanations for this shift cannot be solver independent and cannot center around phenomena observed at the crossover point [41, 43]. In an interesting recent development, Achlioptas, Beame, and Molloy [2] showed that for mixtures of 2-clauses with density  $1 - \epsilon$  and 3-clauses with density 2.28, which are unsatisfiable with high probability, DLL solvers take an exponential time to refute. If  $(1 - \epsilon, 2.28)$  2-clause/3-clause mixtures occur during the solution of a satisfiable 3-SAT instance, then a DLL solver will take time exponential in the order of the instance to solve it. Achlioptas et al. used this to show that a certain DLL solver behaves exponentially at density 3.81. This could also provide an explanation for the observed exponential behavior of GRASP (which is a modified DLL solver) in the low-density region. Cocco and Monasson [14] recently analyzed the computational complexity of random 3-SAT using the  $2 + p$  SAT problem, in which a clause is chosen to be 3-clause with

probability  $p$  and a 2-clause with probability  $1 - p$ . They identified a region to the left of the crossover point, where all instances are almost surely satisfiable, and the complexity of a DLL-based algorithm is exponential. As in [2], it is an appropriate mixture of 2-clauses and 3-clauses, that forces the algorithm to build an exponentially large refutation subtree before finding a solution. Cocco and Monasson show that, depending on the starting density of the 3-SAT instance, a heuristic will or will not avoid building this exponentially large subtree. For low enough density the hard subtree can be avoided with high probability, but at some point it cannot be avoided and we see a transition from polynomial running time to exponential running time. For the GUC heuristic they show that the transition occurs around density 3 (recall that it is known that GUC is linear below 3.003 [11, 23]). All these agree with our observations that the polynomial to exponential behavior occurs in the satisfiable region and is solver-dependent.

### 3. Experimental Setup

Our experimental setup is identical to that of [17, 44]. We generate  $dn$  clauses, each by picking three distinct variables (out of  $n$ ) at random and choosing their polarity uniformly. For each studied point in the  $d \times n$  quadrant we generate at least 100 random instances and apply our solvers. Our experiments were run on Sun Ultra 1 machines. As in [44], we chose to focus on median running time rather than mean running time. The difficulty of completing the runs on very hard instances makes it less practical to measure the mean. Furthermore, the median and the mean are typically quite close to each other, except for the regions that display heavy-tail phenomena, where the median and the mean diverge dramatically [43]. It would be interesting to analyze our data at percentiles other than the 50th percentile (the median) (cf. [43]), though a meaningful analysis for high percentiles would require many more sample points than we have in our experiments.

For the statistical analysis and plotting of data, we used MATLAB, which is an integrated technical computing environment that combines numeric computation, advanced graphics and visualization, and a high-level programming language. The MATLAB ([www.mathworks.com](http://www.mathworks.com)) functions we used for statistical analysis were:

- *polyfit*, for computing the best linear, quadratic, or cubic fit to the data (or the logarithm of the data) using polynomial regression, and
- *corrcoef*, for computing  $r^2$ , the square of correlation ( $r^2$  is the fraction of the variance of one variable that is explained by regression on the other variable) [34].

For each of our data sets we tried to fit:

- a linear curve on the logarithm of both coordinates (notice that this corresponds to a fit of the form  $y = ax^c$  to the data)
- a polynomial curve of the form  $y = ax^{c'} + b$ , where  $c'$  is an integer close to  $c$  of the previous fit, and
- a linear curve on the logarithm of the  $y$ -coordinate, while keeping the  $x$ -coordinate as is (this corresponds to a fit of the form  $y = ae^{cx}$  to the data).

For each fit we computed the  $r^2$  as a measurement of the quality of the fit. Unless stated otherwise, for the results reported in this paper,  $r^2$  exceeded 0.98. This establishes high confidence in the validity of the fit of the curve to the data points.

#### 4. Random 3-SAT and GRASP

GRASP [38] is a SAT solver that augments the basic backtracking search with a conflict-analysis procedure. In order to cut down the search space, a dynamic-learning mechanism based on diagnosing the causes of the conflicts is used. By analyzing conflicts and discovering their causes, GRASP can backtrack non-chronologically to earlier levels in the search tree, potentially pruning large portions of the search space. Moreover, by recording the causes of conflicts, GRASP can avoid running into similar conflicts later during the search.

The experiments described in this section were run on a Sun Ultra 1 with a 167MHz UltraSPARC processor and 128MB RAM. Some changes were made to the default GRASP configuration; we increased the maximum number of backtracks allowed to 1,000,000 and the maximum number of conflicts allowed to 2,000,000. CPU time limit was set to 10,800 seconds. These changes were necessary in order to limit the portion of SAT instances on which GRASP aborted. This artificially lowers our measurements of mean running time, but does not affect our measurements of median running time.

The goal of the experiments was to evaluate GRASP's performance on an initial quadrangle of the  $d \times n$  quadrant. We explored densities from 0.9 to 15. The order of the instances explored depends on the density:

- Density 0.9: 2000 variables (25 variables per step)
- Densities 1, 2, 3, 3.4, and 3.5: 1000 variables (10 variables per step)
- Density 3.6: 800 variables (10 variables per step)
- Density 3.7: 480 (10 variables per step)
- Density 3.8: 450 variables (10 variables per step)
- Density: 4.26: 170 variables (10 variables per step)
- Density 5: 210 variables (10 variables per step)
- Densities 4, 6–15: 250 variables (10 variables per step)

In Figure 1 the median running time is shown on a logarithmic (base 2) scale. (For densities 4.26 and 5 we extrapolated the data up to 250 variables).

We analyzed the median running time as a function of the order for fixed density instances. For low densities (at or below 3.5), our data indicate a quadratic running time. See Figures 1 and 2, where we plot the median running time as a function of the order for instances of density 0.9 and 3.5, respectively. The quadratic behavior of GRASP at low densities should be contrasted with the linear running time at low densities that was reported in [16, 43]. It seems that GRASP's conflict-analysis component has a quadratic overhead.

At densities 3.8 and above, the median running time is exponential in the order, i.e., it behaves as  $2^{\alpha n}$ , where the exponent  $\alpha$  depends on  $d$  (see discussion below). See Figure 2 where we plot the median running time as a function of the order for instances of density

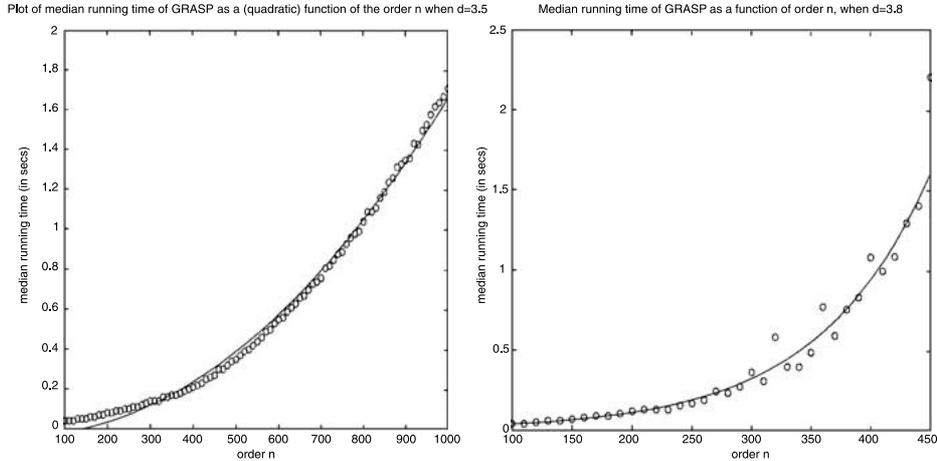


Figure 2. GRASP – median running time for density 3.5 (left) and density 3.8 (right) as a function of the order of the instances. At density 3.5, the best fit curve is quadratic in the order, while at 3.8, the best fit curve is exponential in the order.

3.8.<sup>2</sup> Thus, a phase transition seems to occur between densities 3.5 and 3.8, where the median running time shifts from polynomial to exponential. As the density is increased beyond 3.8, the exponent  $\alpha$  also increases. It peaks at around density 4.26, after which it declines with increased density. Thus, we observe two phase transitions. The second one, in which the exponent reaches its peak, essentially coincides with the crossover point, at which the probability of satisfiability is 0.5. This is the phase transition that was reported at [44] and then studied extensively. This transition, however, is preceded by a more significant one, near density 3.8, where we observe a qualitative shift in the behavior of GRASP. A transition from polynomial to exponential behavior in graph coloring was conjectured in [29] and counter-conjectured in [19]. Such a transition in random 3-SAT near the crossover point is claimed in [16]; this claim, however, was removed in a later paper [17]. We believe that we are the first to demonstrate such a transition in random 3-SAT, and to show that it occurs significantly below the crossover point. The more recent works of Achlioptas et al. [2] and Cocco and Monasson [14] attempt to explain (based on the  $2 + p$ -SAT analysis) why such a transition happens to the left of the crossover point for GRASP (see discussion in Section 2). Comparing with the results in [14], it seems that GRASP is more successful than GUC in avoiding hard subtrees, pushing the transition from polynomial to exponential to a higher density.

The phase transition near 3.8 is accompanied by a “heavy-tail phenomenon”, which is a prevalence of *outliers*, i.e., instances on which the actual running time is at least an order of magnitude (10) larger than the median running time, as well as a divergence of the mean and the median.<sup>3</sup> See Figure 3, where we plot the mean-to-median ratio and the proportion of outliers as a function of the density. The plots show a drastic change in the region between density 3.7 and density 4.3. Both plots show a quick rise and decline. The mean-to-median ratio peaks at around density 4.0 and the proportion of outliers peaks at around density 4.2. For densities between 3.7 and 4.0 we found it quite difficult to analyze the

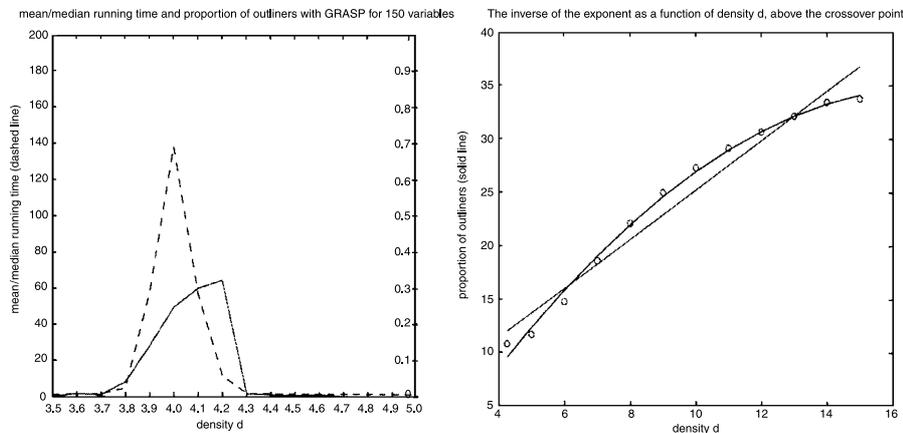


Figure 3. GRASP – (left) Ratio of mean-to-median running time and the proportion of outliers, and (right) the exponent  $1/\alpha$  of median running time as a function of density.

median running time as a polynomial (of low degree) or an exponential function of the order (note the lower  $r^2$  of 0.95 for density 3.8).

Our data suggest that as the density increases from 3.7 to 4.3, random 3-SAT formulas go through a series of changes and perhaps more than one phase transition. The heavy-tail phenomenon for random 3-SAT deserves further study (with many more samples per point in the  $d \times n$  quadrant) to confirm our findings. In particular, the divergence of the two peaks in Figure 3 needs to be reconfirmed or refuted.

As noted above, beyond density 4.26 the exponent  $\alpha$  declines. A theoretical analysis suggests that for DLL solvers  $\alpha$  may decline inversely linearly, i.e., as  $\frac{c}{d}$ , for some constant  $c$ , see [3]. Our data, however, suggest a slower decline, even though one may expect GRASP to be faster than DLL solvers. See Figure 3, where we plot  $\frac{1}{\alpha}$  as a function of  $d$ . Thus, GRASP is not as efficient in the high-density region as it could be. (We should caution, however, that we only have 11 data points, and these data points themselves have been obtained by fitting a linear curve to the logarithm of the median running time. Thus, the finding of a slower decline should be viewed as quite preliminary.)

## 5. Random 3-SAT and CPLEX

The CPLEX MIP Solver is a commercial linear-programming solver for integer variables. It employs a branch-and-bound technique starting from a linear-programming relaxation of the given integer-programming problem. This may be complemented with the dynamic generation of cutting planes [5, 6].

The experiments described in this section were run on a Sun Ultra 1 with a 167MHz UltraSPARC processor and 64 MB RAM. SAT problems were encoded as 0-1 integer-programming problems. Values true and false are represented as 1 and 0. For a clause to be true the sum of the representations of the literals has to be greater or equal to 1. For example, the clause  $\neg x_1 \vee x_2 \vee \neg x_3$  is represented by the inequality  $(1 - x_1) + x_2 + (1 - x_3) \geq 1$ .

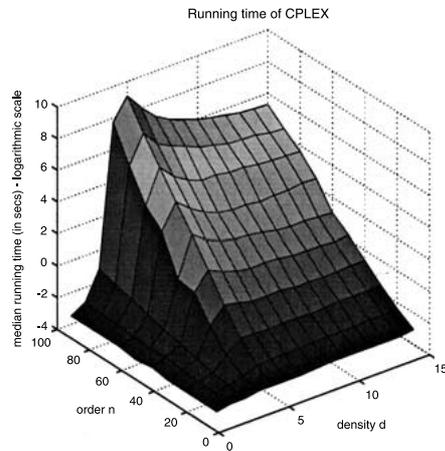


Figure 4. CPLEX – 3-D Plot of median running time.

We used CPLEX to solve problems for densities from 0.9 to 15. The order of the instances was chosen according to the density:

- Densities 0.9, 1.5, 1.6, 1.7 and 1.8: 2000 variables (25 variables per step)
- Density 1: 10000 variables (50 variables per step)
- Density 2: 1800 variables (25 variables per step)
- Density 2.5: 460 variables (10 variables per step)
- Density 3: 250 variables (10 variables per step)
- Density 3.5: 150 variables (10 variables per step)
- Densities 4, 4.26, and 5–15: 120 variables (10 variables per step)

In Figure 4, the median running time is shown on a logarithmic (base 2) scale. Note that the peak at the crossover point is much less pronounced than the one in Figure 1.

The median running time was analyzed as a function of the order for fixed density-instances. For low densities (below 1.7) our data indicate a linear running time. See Figure 5 for median running times for instances of density 1 with up to 10000 variables. For density 2 the median running time is quadratic, while for density 2.5 the running time is cubic. See Figure 5 for median running time for instances of density 2, where for order above 400 the behavior is quadratic. See Figure 6 (left) for median running time for instances of density 2.5, where the behavior is cubic. Thus we seem to have two phase transitions, corresponding to a shift from a linear to quadratic behavior between densities 1 and 2 and a subsequent shift to a cubic behavior between densities 2 and 2.5. The first shift may coincide with the phase transition proved in [8, 40] around density 1.63, as described in Section 2.

At densities 4.0 and above, see Figure 6 (right), the median running time is exponential in the order, i.e., it behaves as  $2^{\alpha n}$ , where the exponent  $\alpha$  depends on  $d$ . As with GRASP, a phase transition seems to occur between densities 2.5 and 4.0. It corresponds to the shift from polynomial to exponential behavior. Again, we believe that as with GRASP, this shift is related to the  $2 + p$ -SAT results in [2, 14]. Note that the polynomial to exponential running time shift for CPLEX occurs in the same region (near density 3.0)

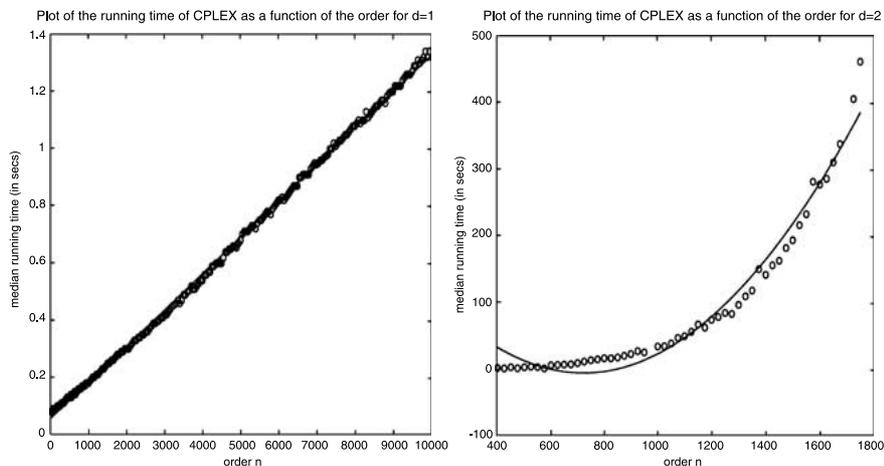


Figure 5. CPLEX – median running time for density 1 (left) and density 2 (right) as a function of the order of the instances.

that the shift for GUC occurs. While CPLEX is a branch-and-bound technique (that can be related to DLL-like heuristics), it also uses cutting-planes technique. Unfortunately, as CPLEX is a commercial tool, we have little access to its underlying algorithms and heuristics, which makes it difficult to offer a precise analysis of its behavior.

As with GRASP, the polynomial-to-exponential transition is accompanied by heavy-tail phenomena. See Figure 7, where we plot the mean-to-median ratio and the proportion of outliers as a function of the density. Note that the heavy-tail phenomenon for CPLEX is not as marked as with GRASP; both peak mean-to-median ratio and peak proportion of

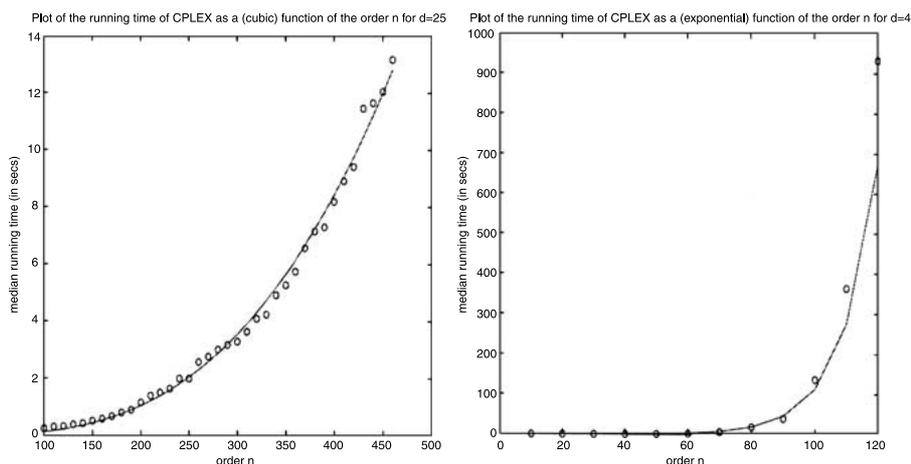


Figure 6. CPLEX – median running time for density 2.5 (left) and density 4 (right) as a cubic and exponential function respectively, of the order of the instances.

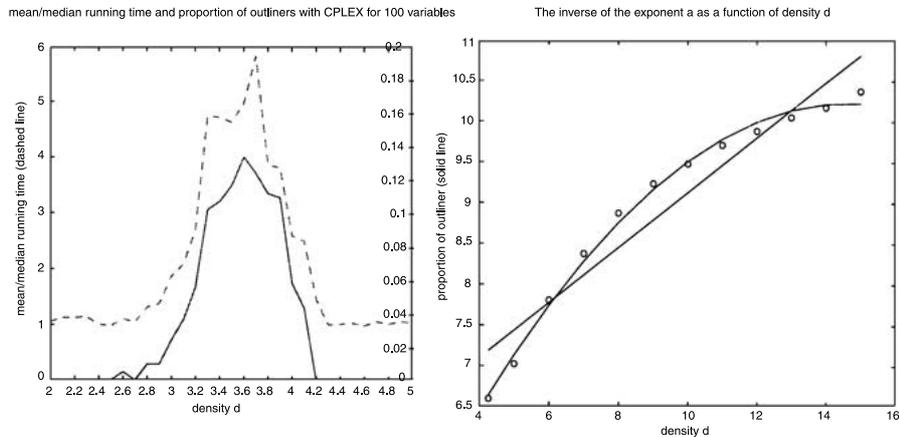


Figure 7. CPLEX – (left) Ratio of mean-to-median running time and proportion of outliers, and (right) the exponent  $1/\alpha$  of median running time as a function of density.

outliers are lower for CPLEX than for GRASP. Note also that the peak for CPLEX occurs at lower densities (around 3.6) than for GRASP (around 4.0).

As with GRASP, the exponent  $\alpha$  peaks at density 4.26 and then declines. Again, our data show a slower decline than  $\frac{c}{d}$ , as suggested in [3] (though the analysis there is for resolution-based procedures, which are weaker than the cutting-planes method used in CPLEX.) See Figure 7, where we plot  $\frac{1}{\alpha}$  as a function of  $d$ .

## 6. Random 3-SAT and CUDD

CUDD [46] is a package that provides functions for the manipulation of Boolean functions, based on the reduced, ordered, binary decision diagram (ROBDD) representation [9]. A binary decision diagram (BDD) is a rooted directed acyclic graph that has only two terminal nodes labeled 0 and 1. Every non-terminal node is labeled with a Boolean variable and has two outgoing edges labeled 0 and 1. An ordered binary decision diagram (OBDD) is a BDD with the constraint that the input variables are ordered and every path in the OBDD visits the variables in ascending order. An ROBDD is an OBDD where every node represents a distinct logic function.

Unlike GRASP and CPLEX, CUDD does not search for a satisfying truth assignment. Rather, it constructs a compact symbolic representation of the set of *all* satisfying truth assignments. Then, the resulting ROBDD is compared against the predefined constant 0 in order to find if an instance is (un)satisfiable. It is important to note that very large sets of truth assignments can have very compact ROBDD representation [9], which explains the effectiveness of ROBDDs in hardware verification [10, 31]. As we see later, CUDD performs well in the very-low-density region, where the set of satisfying truth assignments is very large.

The experiments described in this section were run on a Sun Ultra 1 with a 167MHZ UltraSPARC processor and 64MB RAM. The CUDD package has been used through the GLU C-interface [48], a set of low-level utilities to access BDD packages. It is well known

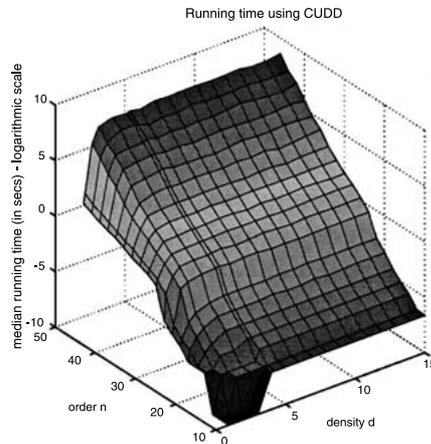


Figure 8. CUDD – 3-D Plot of median running time.

that the size of the ROBDD for a given function depends on the variable order chosen for that function. We have used automatic dynamic reordering during the tests with the default method for automatic reordering of CUDD.

As in the preceding two sections, the goal of the experiments was to evaluate CUDD's performance on an initial quadrangle of the  $d \times n$  quadrant. We explored densities 0.1, 0.5, and 1 to 15. The order of the instances explored depends on the density:

- Density 0.1: 1480 variables (10 variables per step)
- Density 0.5: 136 variables (2 variables per step)
- Density 0.9 and 1: 68 variables (2 variables per step)
- Densities 1.5, 2–4, 4.26, 5–15: 46 variables (2 variables per step)

In Figure 8 the median running time is shown on a logarithmic (base 2) scale. Note the absence of a peak (contrast with Figures 1 and 4).

We analyzed the median running time as a function of the order for fixed-density instances. At densities 0.5 and above, the median running time is exponential in the order, i.e., it behaves as  $2^{\alpha n}$ . See Figure 9 (right) for median running time for instances of density 1, where the behavior is exponential. At density 2 and above the exponent  $\alpha$  is independent of the density. In particular, there seems to be nothing special about the crossover point at density 4.26. The explanation for this behavior is that the running time of ROBDD-based algorithms is determined mostly by the size of the manipulated ROBDDs. Our algorithm involves  $dn$  product operations between a possibly large ROBDD (representing all truth assignments of the clauses processed so far) and a small ROBDD (representing seven truth assignments of the currently processed clause). Thus, the running time of our algorithm is determined by the largest intermediate ROBDD constructed. As is shown in Figure 10, the peak in ROBDD size is attained after processing about  $2n$  clauses, which explains the flattening of the running-time plot at density 2, and suggests that a phase transition in terms of ROBDD size occurs at about this density.

As ROBDDs are symmetrical with respect to the set they represent and its complement, both very small sets and very large sets can be represented by small ROBDDs [9].

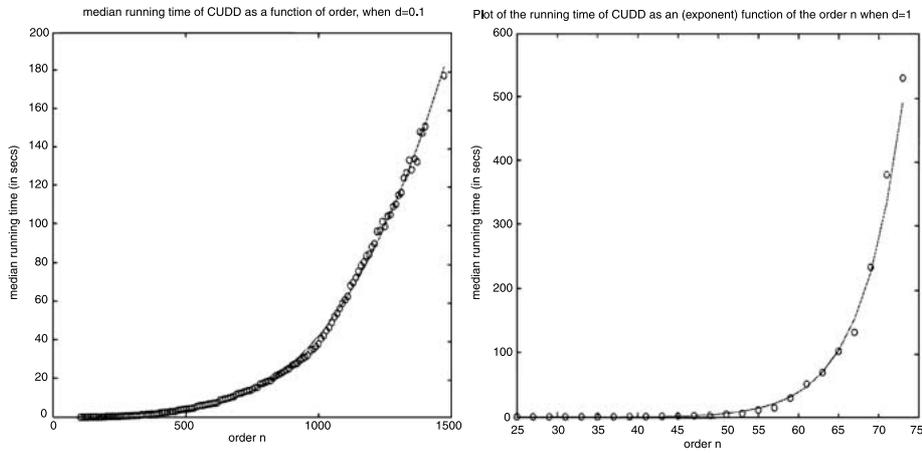


Figure 9. CUDD – median running time for density 0.1 (left) and for density 1 (right) as a cubic and an exponential respectively function of the order of the instances.

This suggests that we may see polynomial behavior for very low density instances, which have a large number of satisfying truth assignments. To check this conjecture we measured the median running time of CUDD for instances of density 0.1. Our results indicate a cubic-time behavior, see Figure 9. This suggests the existence of another phase transition between densities 0.1 and 0.5. This result should be contrasted with that of [45], in which the running time for explicitly enumerating all solutions of random constraint-satisfaction instances increases as the density decreases.

Unlike GRASP and CUDD, we did not observe a heavy-tail phenomenon with CUDD: there are no outliers and the mean-to-median ratio is independent of the density (see Figure 10).

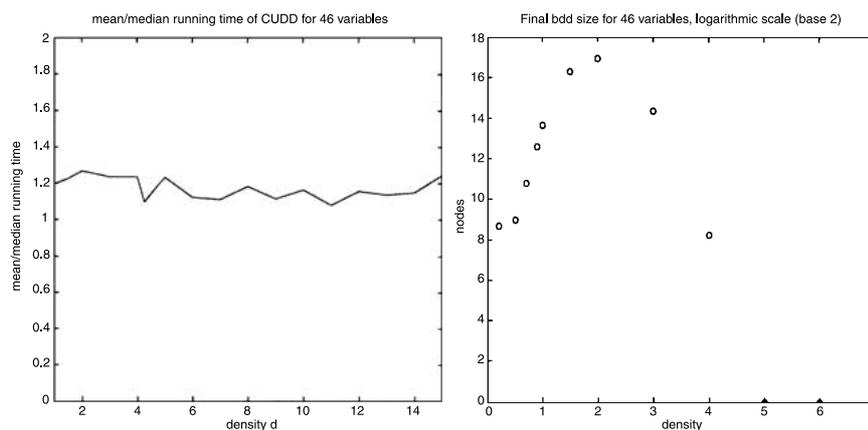


Figure 10. CUDD – (left) Ratio of mean-to-median running time and (right) median ROBDD size as a function of density.

## 7. Discussion

We provide experimental evidence for the following hypotheses. First, the connection between the phase transition in computational complexity and the phase transition in satisfiability is not as tight as has been claimed. It is not the case that the shift from polynomial to exponential complexity occurs at or very close to the crossover point, as has been widely believed [43, 44]. Second, not only does the density at which the shift from polynomial to exponential time complexity vary with the choice of solver, but the very shape of the surface of the median running time (an experimental surrogate for average-time complexity), as a function of the density  $d$  and the order  $n$ , changes with the solver. Finally, the density-order quadrant contains several phase transitions; in fact, the region between density 0 and density 4.26 seems to be rife with phase transitions, which are also solver dependent. In essence, each solver provides us with a different tool with which to study the complexity of random 3-SAT. This is analogous to astronomers observing the sky using telescopes that operate at different wave lengths. We thus hope to alleviate the “fixation” with DLL solvers and the crossover point at 4.26.

Our experiments reveal a marked difference between solvers like GRASP and CPLEX which are search based, and GRASP and CPLEX display interesting similarities in the shapes of the median running time surface despite their algorithmic differences, and ROBDD-based solvers, like CUDD, which are based on compactly representing all satisfying truth assignments. While the interesting region for GRASP and CPLEX is between 2.5 and 4.3, the interesting region for CUDD occurs below density 2. This refutes earlier conjectures (cf. [36]) that the peak in median running time around the crossover point is essentially solver independent. For both GRASP and CPLEX, we observed a new phase transition where the median running time shifts from being polynomial in the order to being exponential in the order. For GRASP the transition occurs at around density 3.8, while for CPLEX the transition occurs earlier, near density 3.0. From the perspective of average-time complexity this is a significant phase transition because it corresponds to a qualitative shift in the behavior of the solver. We also observed several other phase transitions for CPLEX and for CUDD. This suggests that it would be interesting to explore the behavior of other SAT solvers, such as RELSAT [32] or SATZ [37], on the  $d \times n$  quadrant.

With fine grained sampling of the density parameter, and by exploring a greater range in the number of variables, we can start to document for each solver, phase transitions that correspond to significant shifts in the shape of the running time of the solver. These phase transitions are important to our understanding of the computational complexity of random 3-SAT, and can be used as a justification to develop density-based solvers for 3-SAT, i.e., solvers which use information about the density of an instance, to choose the most appropriate algorithmic technique. (Of course, the applicability of average-case complexity results to a real-world setting depends on the fit of the assumed instance distribution to the application at hand. From that perspective, the fixed-width distribution studied here is only one of many possible distributions.)

While our results are purely empirical, as the lack of success with formally proving a sharp complexity threshold at the crossover point indicates (cf. [1, 21, 22]), providing

rigorous proof for our qualitative observations may be a very difficult task, especially for sophisticated solvers like the ones studied in this paper.

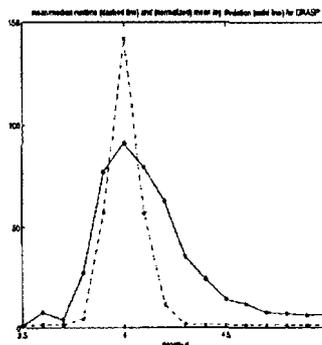
### Acknowledgments

The last author is grateful to Phokion G. Kolaitis for stimulating discussions.

Demopoulos and Vardi were supported in part by NSF grants CCR-9700061 and IIS-9908435, and by a grant from the Intel Corporation. San Miguel Aguirre's current address is: Dept. of Computer Science, Instituto Tecnológico Autónomo de México, Río Hondo 1, 01000 Mexico City. He did this work while on sabbatical at Rice University, supported in part by CONACyT grant 145502. Subramanian was supported in part by NSF grant IRI-9796046.

### Notes

1. It is easy to see that 3-SAT is NP-complete for instances of each fixed density, as the generic reduction of NP to 3-SAT [24] produces instances of fixed density and each density can be obtained by adding linearly many redundant variables or redundant clauses. Thus, we'd expect the worst-cases running time to be exponential for all densities, and, consequently, to find hard instances in the "easy" region, cf. [25]. The issue of median vs. mean running time is discussed later.
2. The  $r^2$  for this plot is 0.95, while the  $r^2$  for all of the polynomial fits that we tried was  $\leq 0.9$ . That gives us confidence that the running time is exponential in the order.
3. There are several ways to analyze inequality, asymmetry and outliers among data. One way to measure inequality within a sample set is to use the mean-log deviation  $GE(0)$  [49], which has been used extensively to measure inequality in the context of economic studies of populations and income distributions. It is defined as  $GE(0) = \frac{1}{n} \sum_{i=1}^n \log\left(\frac{y_i}{\bar{y}}\right)$ , where  $n$  is the number of individuals in the sample,  $y_i$  is the income of individual  $i$  and  $\bar{y}$  the mean income. See the figure below where we plotted the mean-log deviation for the running time of GRASP against the mean-over-median runtime ratio, in the density region of 3.5 to 5. Mean-log deviation shows a behavior similar to the mean-over-median ratio, and also with the number of outliers (see Figure 3). It peaks at density 4 and it indicates a dramatic increase of inequality to the left of the crossover point. Also, kurtosis and skewness [34] can be used to show whether the data are peaked or flat relative to a normal distribution and whether data are symmetrical or skewed to the right or left. Calculations of kurtosis and skewness on the running time data of GRASP show them both to rise quickly and then quickly drop again; they both peak at around density 3.6. These indicate, as the rest of the statistics reported in this paper, that a significant number of outliers appear between densities 3.5 and 4.



## References

1. Achlioptas, D. (2000). Setting two variables at a time yields a new lower bound for random 3-SAT. In *Proc. 32th ACM Symp. on Theory of Computing*, pages 28–37.
2. Achlioptas, D., Bename, P., & Molloy, M. (2001). A sharp threshold in proof complexity. In *Proc. 33rd ACM Symp. on Theory of Computing*, pages 337–346.
3. Beame, P., Karp, R. M., Pitassi, T., & Saks, M. E. (1998). On the complexity of unsatisfiability proofs for random  $k$ -CNF formulas. In *Proc. 30th ACM Symp. on Theory of Computing*, pages 561–571.
4. Biere, A., Cimatti, A., Clarke, E. M., Fujita, M., & Zhu, Y. (1999). Symbolic model checking using SAT procedures instead of BDDs. In *Proc. 36th Conf. on Design Automation*, pages 317–320.
5. Bixby, R. E. (1992). Implementing the Simplex method: the initial basis. *ORSA J. on Computing*, 4(3): 267–284.
6. Bixby, R. E. (1994). Progress in linear programming. *ORSA J. on Computing*, 6(1): 15–22.
7. Bouquet, F. (1999). *Gestion de la dynamique et énumération d'implicants premiers: une approche fondée sur les Diagrammes de Décision Binaire*. PhD thesis, Université de Provence, France.
8. Broder, A. Z., Frieze, A. M., & Upfal, E. (1993). On the satisfiability and maximum satisfiability of random 3-CNF formulas. In *Proc. 4th Annual ACM-SIAM Symp. on Discrete Algorithms*, pages 322–330.
9. Bryant, R. E. (1986). Graph-based algorithms for Boolean function manipulation. *IEEE Trans. on Computers*, 35(8): 677–691.
10. Burch, J. R., Clarke, E. M., McMillan, K. L., Dill, D. L., & Hwang, L. J. (1992). Symbolic model checking:  $10^{20}$  states and beyond. *Information and Computation*, 98(2): 142–170 (June).
11. Chao, M., & Franco, J. V. (1990). Probabilistic analysis of a generation of the unit-clause literal selection heuristics for the  $k$ -satisfiability problem. *Information Sciences*, 51: 289–314.
12. Cheeseman, P., Kanefsky, B., & Taylor, W. M. (1991). Where the really hard problems are. In *Proc. 12th Int'l Joint Conf. on Artificial Intelligence (IJCAI '91)*, pages 331–337.
13. Chvátal, V., & Szemerédi, E. (1988). Many hard examples for resolution. *J. of the ACM*, 35(4): 759–768.
14. Cocco, S., & Monasson, R. (2001). Statistical physics analysis of the computational complexity of solving random satisfiability problems using backtrack algorithms. *European Physical Journal B*, 22: 505–531.
15. Cook, S. A., & Mitchell, D. G. (1997). Finding hard instances of the satisfiability problem: a survey. In Du, Gu, & Pardalos, eds., *Satisfiability Problem: Theory and Applications*, Volume 35 of *Dimacs Series in Discrete Math. and Theoretical Computer Science*. AMS.
16. Crawford, J. M., & Auton, L. D. (1993). Experimental results on the crossover point in satisfiability problems. *AAAI*, pages 21–27.
17. Crawford, J. M., & Auton, L. D. (1996). Experimental results on the crossover point in random 3-SAT. *Artificial Intelligence*, 81(1–2): 31–57.
18. Crawford, J. M., & Baker, A. B. (1994). Experimental results on the application of satisfiability algorithms to scheduling problems. *AAAI*, 2: 1092–1097.
19. Cullberson, J., & Gent, I. P., (2000). Frozen development in graph coloring. Technical report, Department of Computer Science, University of St. Andrews. APES Research Report APES-19-20004.
20. Davis, M., Logemann, G., & Loveland, D. (1962). A machine program for theorem proving. *Comm. of the ACM*, 5: 394–397.
21. Dubois, O., Boufkhad, Y., & Mandler, J. (2000). Typical random 3-SAT formulae and the satisfiability threshold. In *Proc. 11th Annual ACM-SIAM Symp. on Discrete Algorithms*, pages 126–127.
22. Friedgut, E. (1999). Necessary and sufficient conditions for sharp threshold of graph properties and the  $k$ -SAT problem. *J. Amer. Math. Soc.*, 12: 1917–1054.
23. Frieze, A., & Suen, S. (1996). Analysis of two simple heuristics for random instances of  $k$ -SAT. *J. of Algorithms*, 20(2): 312–355.
24. Garey, M. R., & Johnson, D. S. (1979). *Computers and Intractability, A Guide to the Theory of NP-Completeness*. New York, NY: W. H. Freeman.
25. Gent, I. P., & Walsh, T. (1994). Easy problems are sometimes hard. *Artificial Intelligence*, 70(1–2): 335–345.
26. Gent, I. P., & Walsh, T. (1994). The SAT phase transition. In *Proc. European Conf. on Artificial Intelligence*, pages 105–109.
27. Groote, J. F. (1996). Hiding propositional constants in BDDs. *Formal Methods in System Design*, 8: 91–96.

28. Groote, J. F., & Zantema, H. (2000). Resolution and binary decision diagrams cannot simulate each other polynomially. Technical report, Department of Computer Science, Utrecht University. Technical Report UU-CS-2000-14.
29. Hogg, T., & Williams, C. P. (1994). The hardest constraint problems: a double phase transition. *Artificial Intelligence*, 69(1–2): 359–377.
30. Hooker, J. N. (1988). Resolution vs. cutting plane solution of inference problems: some computational experience. *Operations Research Letters*, 7: 1–7.
31. Jha, S., Lu, Y., Minea, M., & Clarke, E. M. (1997). Equivalence checking using abstract BDDs. In *Proc. 1997 Int'l Conf. on Computer Design*, pages 332–337.
32. Bayardo, Roberto J. Jr., & Schrag, Robert. (1997). Using CSP look-back techniques to solve real-world SAT instances. In *AAAI/IAAI*, pages 203–208.
33. Kautz, H., & Selman, B. (1992). Planning as satisfiability. In *Proc. European Conference on Artificial Intelligence*, pages 359–379.
34. Kirk, R. E. (1999). *Statistics: An introduction*, fourth edition. Harcourt Brace, Fort Worth.
35. Kirkpatrick, S., & Selman, B. (1994). Critical behavior in the satisfiability of random formulas. *Science*, 264: 1297–1301.
36. Larrabee, T., & Tsuji, Y. (1992). Evidence for a satisfiability threshold for random 3CNF formulas. Technical report, University of California, Santa Cruz. Technical report USCS-CRL-92042.
37. Chu Min Li, & Anbulagan. (1997). Heuristics based on unit propagation for satisfiability problems. In *IJCAI (I)*, pages 366–371.
38. Marques Silva, J. P., & Sakallah, K. A. (1999). GRASP-A search algorithm for propositional satisfiability. *IEEE Trans. on Computers*, 48(5): 506–521.
39. Mitchell, D. G., & Levesque, H. J. (1996). Some pitfalls for experimenters with random SAT. *Artificial Intelligence*, 81(1–2): 111–125.
40. Mitzenmacher, M. (1997). Tight thresholds for the pure literal rule. Technical report, Digital System Research Center, Palo Alto, California. SRC Technical Note 1997 - 011.
41. Monasson, R., Zecchina, R., Kirkpatrick, S., Selman B., & Troyansky, L. (1999). Determining computational complexity from characteristic ‘phase transitions’. *Nature*, 400: 133–137.
42. Nam, G. J., Sakallah, K. A., & Rutenbar, R. A. (1999). Satisfiability-based layout revisited: Detailed routing of complex FPGAs via search-based boolean sat. In *Proc. ACM Int'l Symp. on Field-Programmable Gate Arrays (FPGA'99)*, pages 167–175.
43. Selman, B., & Kirkpatrick, S. (1996). Critical behavior in the computational cost of satisfiability testing. *Artificial Intelligence*, 81(1–2): 273–295.
44. Selman, B., Mitchell, D. G., & Levesque, H. J. (1996). Generating hard satisfiability problems. *Artificial Intelligence*, 81(1–2): 17–29.
45. Smith, B. M., & Dyer, M. E. (1996). Locating the phase transition in binary constraint satisfaction problems. *Artificial Intelligence Journal*, 8(1–2): 155–181.
46. Somenzi, F. (1998). CUDD: CU Decision Diagram package. release 2.3.0. Dept. of Electrical and Computer Engineering. University of Colorado at Boulder.
47. Svenson, P., & Nordahl, M. G. (1999). Relaxation in graph coloring and satisfiability problems. *Phys. Rev. E*, 59(4): 3983–3999.
48. The VIS Group. (1996). VIS: a system for verification and synthesis. In Alur, R., & Henzinger, T., eds., *Proc. 8th Int'l Conf. on Computer Aided Verification (CAV '96)*, LNCS 1102, pages 428–432.
49. Theil, H. (1979). The measurement of inequality by component of income. *Economics Letter*, 2.
50. Uribe, T. E., & Stickel, M. E. (1994). Ordered binary decision diagrams and the Davis-Putnam procedure. In *First Int'l Conf. on Constraints in Computational Logics*, Vol. 845 of *Lecture Notes in Computer Science*, pages 34–39, Munich: Springer-Verlag (September).