# The Resolution Complexity of Random Constraint Satisfaction Problems

Michael Molloy
Department of Computer Science
University of Toronto
molloy@cs.toronto.edu

Mohammad R. Salavatipour
Department of Computing Science
University of Alberta
mreza@cs.ualberta.ca

September 1, 2006

#### Abstract

We consider random instances of constraint satisfaction problems where each variable has domain size d, and each constraint contains t restrictions on k variables. For each (d, k, t) we determine whether the resolution complexity is a.s. constant, polynomial or exponential in the number of variables. For a particular range of (d, k, t), we determine a sharp threshold for resolution complexity where the resolution complexity drops from a.s. exponential to a.s. polynomial when the clause density passes a specific value.

#### 1 Introduction

A constraint satisfaction problem (CSP) is a generalized form of satisfiability which is widely studied in the Artificial Intelligence community. For example, the journal *Constraints* is devoted to these problems. Roughly speaking, a CSP generalizes SAT in the sense that variables can draw their values from a more general domain than simply  $\{T, F\}$ , and each clause (a.k.a. constraint) forms a set of restrictions on the values that the variables in the clause may jointly take.

Random instances of k-SAT have been extremely well-studied over the past few decades (see [1] for many references). More recently, the interest in this area has expanded into random instances of various generalizations of k-SAT, such as NAE-SAT[3], XOR-SAT[14, 18, 19], (2 + p)-SAT [36, 38, 39, 4, 2] and many others. All of these can be expressed as CSP's. It was natural for this interest to eventually spread to random instances of CSP's, rigorously in [5, 15, 34, 20, 32, 33, 41] and experimentally even earlier (see [24] for a good survey).

One of the most important results regarding random k-SAT is that of Chvátal and Szemerédi[13], who showed that for any  $k \geq 3$  and c > 0, a random instance of k-SAT with n variables and cn clauses will almost surely (a.s.)<sup>1</sup> have no resolution proof of unsatisfiabilty of length less than  $2^{\Theta(n)}$ . It is easy to show that for large values of c, such random instances are almost surely unsatisfiable. This immediately implied that for sufficiently large values of c, any Davis-Putnam style algorithm will take exponential time on such an input. Furthermore, it provided an astoundingly vast and rich class to the, beforehand rather sparse[26, 40], list of unsatisfiable instances of k-SAT for which there is no polytime resolution proof of unsatisfiability. Such instances are of great interest since their existence can be viewed as a step toward proving that there are some unsatisfiable instances

<sup>&</sup>lt;sup>1</sup>Formal definitions of these and other terms will appear in the next section

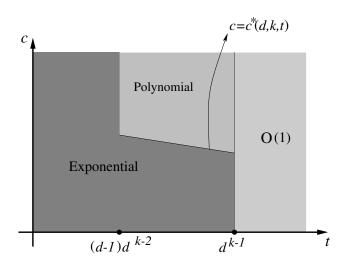


Figure 1: Resolution complexities

with no polytime proof of unsatisfiability of any kind, i.e. that  $NP \neq co - NP$ . Chvátal and Szemerédi's paper spawned numerous extensions and generalizations, eg. [7, 8, 2, 6], including a general framework for proving lower bounds on resolution complexity by Ben-Sasson and Wigderson[9].

Mitchell [32, 33] extended the framework of Ben-Sasson and Wigderson to the setting of CSP's. He then used this framework to prove exponential lower bounds on the resolution complexity of a very natural class of random CSP's - one where the number of restrictions per constraint is fixed. Specifically, he considered random CSP's with domain size  $d \geq 2$ , and every constraint containing precisely t restrictions on  $k \geq 2$  variables.<sup>2</sup> Note that these CSP's are trivial if either d or k is equal to 1, and that they are the well-studied 2-SAT when d = k = 2 and t = 1. Mitchell showed that for  $t \leq (d-1)/2$ , k = 2 and for  $t \leq d-1$ ,  $k \geq 3$ , and for any constant c > 0, such a random instance with cn constraints will almost surely have no subexponential proof of unsatisfiability.<sup>3</sup> Again, it is easy to see that for sufficiently large c, these instances are a.s. unsatisfiable (for t > 0). In contrast, Achlioptas et al[5] showed that for  $t \geq d^{k-1}$ , and any c > 0, such a random instance will a.s. have an unsatisfiable subproblem of size O(1), and thus will have a O(1)-length resolution proof of unsatisfiability. In this paper, we fill in the gap between d-1 and  $d^{k-1}$ . Using  $\mathcal{F}_{n,M}^{d,k,t}$  to denote such a random CSP with M constraints, we prove the following theorems which are summarized in Figure 1.

**Theorem 1** For any constants  $d, k \geq 2$  and  $1 \leq t < (d-1)d^{k-2}$ , and for every constant c > 0,  $\mathcal{F}_{n,M=cn}^{d,k,t}$  a.s. has resolution complexity at least  $2^{\Theta(n)}$ .

For  $d, k \geq 2$  and  $t \geq (d-1)d^{k-2}$ , we define

$$c^*(d,k,t) = \frac{1}{dk(k-1)} {d^k \choose t} / {d^k - (d-1)d^{k-2} \choose t - (d-1)d^{k-2}}.$$

<sup>&</sup>lt;sup>2</sup>This natural model was, historically, one of the first two random models of random CSP to be studied; the other turned out to be problematic and a.s. has O(1) length resolution proofs of unsatisfiablity for any non-trivial number of constraints. See [34] or [24] for more details; the latter reference contains more than 30 references to the study of the model considered here.

<sup>&</sup>lt;sup>3</sup>In [32], Mitchell claims to prove that this holds for  $t \leq (d-1)(k-1)$  so long as d, k are not both 2. But there is an unfortunate error in his Lemmas 8 and 10, and his proof only holds for  $t \leq (d-1)/2, k=2$  and  $t \leq d-1, k \geq 3$ .

(An explanation of the derivation of this expression will have to wait until the end of the next section - see Lemma 9.)

**Theorem 2** For any constants  $d, k \geq 2$  and  $(d-1)d^{k-2} \leq t < d^{k-1}$ , and for every  $c < c^*(d,k,t)$ ,  $\mathcal{F}_{n,M=cn}^{d,k,t}$  a.s. has resolution complexity at least  $2^{\Theta(n)}$ .

**Theorem 3** For any constants  $d, k \geq 2$  and  $(d-1)d^{k-2} \leq t < d^{k-1}$ , and for every constant  $c > c^*(d, k, t)$ ,  $\mathcal{F}_{n, M = cn}^{d, k, t}$  a.s. has resolution complexity poly(n).

Trivially, the resolution complexity of a satisfiable CSP is infinite, so Theorem 2 is of no interest if for all  $c < c^*(d, k, t)$ ,  $\mathcal{F}_{n,M}^{d,k,t}$  is a.s. satisfiable. This is, in fact, well-known to be the case for d = 2, k = 2, t = 1 (i.e. 2-SAT), and we prove here that it is also the case for d = 2, k = 3, t = 3 (see Theorem 5 below). We prove that it is not the case for d = 2, k = 3, t = 2 and for all other d, k (note that when d = 2, k = 3 and t is as in Theorem 2, then  $t \in \{2, 3\}$ ):

**Theorem 4** (a) For any constants  $d, k \geq 2$  and for every  $c > \ln d / \ln[d^k/(d^k - t)]$ ,  $\mathcal{F}_{n, M = cn}^{d, k, t}$  is a.s. unsatisfiable.

(b) For any constants  $d, k \geq 2, (d, k) \notin \{(2, 2), (2, 3)\}$  and  $(d - 1)d^{k-2} \leq t < d^{k-1}$ :

$$\frac{\ln d}{\ln[d^k/(d^k - t)]} < c^*(d, k, t).$$

(c) For every c > 2.114,  $\mathcal{F}_{n,M=cn}^{2,3,2}$  is a.s. unsatisfiable.

Parts (a,b) of Theorem 4 prove that  $\mathcal{F}_{n,M}^{d,k,t}$  is a.s. unsatisfiable for some values of  $c < c^*(d,k,t)$  for  $d,k \geq 2, (d,k) \notin \{(2,2),(2,3)\}$  and  $(d-1)d^{k-2} \leq t < d^{k-1}$ . Part (c) proves the same for the case d=2,k=3,t=2 since  $c^*(2,3,2)=7/3>2.114$ . The next theorem shows that this is not the case for d=2,k=3,t=3, and that here, just as in 2-SAT, there are short resolution proofs of unsatisfiability for every value of c above the threshold of satisfiability.

**Theorem 5** For every  $c < 7/9 = c^*(2,3,3)$ ,  $\mathcal{F}_{n,M=cn}^{2,3,3}$  is a.s. satisfiable. Thus, 7/9 is the (sharp) threshold of satisfiability for  $\mathcal{F}_{n,M=cn}^{2,3,3}$ .

As mentioned above, when  $t \geq d^{k-1}$ , the resolution complexity of  $\mathcal{F}_{n,M=cn}^{d,k,t}$  is a.s. O(1). So these theorems completely characterize the resolution complexity of  $\mathcal{F}_{n,M=cn}^{d,k,t}$  for every constant d,k,t,c except for  $(d-1)d^{k-2} \leq t < d^{k-1}$  and  $c = c^*(d,k,t)$ . Here we have a sharp threshold for resolution complexity, similar to that found in [2], where the main technical result was:

**Theorem 6** For any  $\Delta, \epsilon > 0$ , consider a random CNF-formula F on n variables with  $\Delta n$  3-clauses and  $(1 - \epsilon)n$  2-clauses where every such formula is equally likely. F a.s. has resolution complexity at least  $2^{\Theta(n)}$ .

(The other side of the "sharp threshold", i.e. that if the number of 2-clauses is  $(1+\epsilon)n$  for some  $\epsilon > 0$  then a.s. the resolution complexity of F is poly(n), was previously known to follow from the work in [12].)

Theorems 2 and 3 are at heart very similar to Theorem 6. For  $(d-1)d^{k-2} \leq t < d^{k-1}$ , a certain type of constraint called a *forcer* arises. Forcers play, essentially, the same role that 2-clauses play in random CNF-formulas. We show that if  $c > c^*$  then the forcers alone provide a unsatisfiable CSP with low resolution complexity, while if  $c < c^*$  then, even along with the additional non-forcer constraints, the CSP has high resolution complexity.

Independently, Gao and Culberson [23] proved Theorem 3 for the special case d=2. Essentially, they showed that in this case the forcers imply 2-clauses and for  $c>c^*$  these 2-clauses form a random instance of 2-SAT which is above the satisfiability threshold. It is well-known that such an instance will have low resolution complexity. We remark more on this at the end of Subsection 2.3.

At first, we tried to adapt the lengthy proof of Theorem 6 to the setting of Theorem 2, but we were unsuccessful. Fortunately, we found an alternate proof technique, and to our pleasant surprise, it produced a proof of Theorem 2 which was dramatically shorter than the proof from [2] of Theorem 6. In fact, our technique yields a short proof of Theorem 6 which we provide in an appendix. This technique looks like it will be of value to those who wish to prove future similar theorems.

We close this section by mentioning that the class of models of random CSP's considered here is a subset of the more general class introduced in [15, 34]. That class contains a much wider range of problems, including XOR-SAT and d-colourability. It would be very nice to characterize which models from that larger class exhibit high resolution complexity, but thus far we are unable to do so.

#### 2 Preliminaries

Here we give formal definitions of some of the concepts discussed in the introduction, along with other concepts required for the remainder of the paper.

#### 2.1 The random model

In our setting, the variables of our problem all have the same domain of permissible values,  $\mathcal{D} = \{1,...,d\}$ , and all constraints will be on k variables, for some fixed integers  $d,k \geq 2$ . Given a k-tuple of variables,  $(x_1,...,x_k)$ , a restriction on  $(x_1,...,x_k)$  is a k-tuple of values  $R = (\delta_1,...\delta_k)$  where each  $\delta_i \in \mathcal{D}$ . A set of restrictions on a k-tuple  $(x_1,...,x_k)$  is called a constraint (or a clause). An assignment of values to the variables of a constraint C satisfies C if that assignment is not one of the restrictions in C. A constraint satisfaction problem (CSP) consists of a set of variables and a collection of constraints on subsets of those variables. An assignment of values to all variables in a CSP satisfies that CSP if every constraint is simultaneously satisfied. A CSP is satisfiable if there is at least one such satisfying assignment. The degree of a variable is the number of constraints that it lies in.

A subproblem of a CSP  $\mathcal{I}$  is a CSP which is obtained by removing some of the variables and some of the constraints from  $\mathcal{I}$ , where of course, if a variable x is removed then every constraint containing x is also removed. When there is no possibility of confusion we often use, for example,  $\mathcal{I} - \{C_1, C_2\}$  to denote the subproblem obtained by deleting the constraints  $C_1, C_2$  from  $\mathcal{I}$  and  $\mathcal{I} - \{x_1, x_2\}$  to denote the subproblem obtained by deleting the variables  $x_1, x_2$  from  $\mathcal{I}$ , along with any constraints containing them.

Recall that a k-uniform hypergraph is a generalization of a graph, where each edge contains k vertices. The constraint hypergraph of a CSP is the k-uniform hypergraph whose vertices correspond to the variables, and whose edges correspond to the k-tuples of variables which have (non-empty)

constraints. Of course, when k = 2, the constraint hypergraph is simply a graph, and so we often call it the constraint graph.

We define  $\Omega^{d,k,t}$  to be the set of CSP's in which every variable has domain  $\{1,...,d\}$ , every constraint has k variables and t restrictions, and no two constraints use the same k-tuple of variables.

**The Random Model:** Specify c, n, d, k, t and let M = cn. First choose a random constraint hypergraph with n vertices and M edges of size k, where each such hypergraph is equally likely. Next, for each edge e, we choose a random constraint on the k variables of e, with domains  $\mathcal{D} = \{1, ..., d\}$ , uniformly from amongst all constraints with exactly t restrictions.

Note that every member of  $\Omega^{d,k,t}$  with n variables and M clauses is equally likely to be chosen. We use  $\mathcal{F}_{n,M}^{d,k,t}$  to denote a random CSP drawn from this model. We say that a property holds almost surely (a.s.) if the probability that it holds tends to 1 as n tends to infinity.

**Remark:** Alternatively, we could have chosen the constraint hypergraph by making an independent choice for each potential edge, deciding to put it in the hypergraph with probability  $p = \frac{c \times k!}{n^{k-1}}$ . We denote the resulting random CSP by  $\mathcal{F}_{n,p}^{d,k,t}$ . This model is, in many senses, equivalent to  $\mathcal{F}_{n,M}^{d,k,t}$ , as we describe in Appendix B. In particular, Lemma 23 implies that all the theorems in this paper translate to  $\mathcal{F}_{n,p}^{d,k,t}$ . We will make use of the equivalence of these models in the proofs of Lemmas 15 and 16.

#### 2.2 Resolution complexity

For a boolean CNF-formula F, a resolution refutation of F with length r is a sequence of clauses  $C_1, ..., C_r = \emptyset$  such that each  $C_i$  is either a clause of F, or is derived from two earlier clauses  $C_j, C_{j'}$  for j, j' < i by the following rule:  $C_j = (A \vee x), C_{j'} = (B \vee \overline{x})$  and  $C_i = (A \vee B)$ , for some variable x. The resolution complexity of F, denoted  $\mathbf{RES}(F)$ , is the length of the shortest resolution refutation of F. (If F is satisfiable then  $\mathbf{RES}(F) = \infty$ .)

Mitchell[33] discusses two natural ways to extend the notion of resolution complexity to the setting of a CSP. These two measures of resolution complexity are denoted  $\mathbf{C} - \mathbf{RES}$  and  $\mathbf{NG} - \mathbf{RES}$ . The latter appears on the surface to be the most natural extension in that it extends resolution rules to the setting of a CSP and then carries them out.  $\mathbf{C} - \mathbf{RES}$ , on the other hand, converts a CSP to a boolean CNF-formula and then carries out CNF-resolution on that formula. Mitchell shows that for every CSP instance  $\mathcal{I}$ ,  $\mathbf{C} - \mathbf{RES}(\mathcal{I}) \leq \text{poly}(\mathbf{NG} - \mathbf{RES}(\mathcal{I}))$  whereas there are many choices for  $\mathcal{I}$  for which the converse it not true. Furthermore, all commonly used resolution-type CSP algorithms correspond nicely to the  $\mathbf{C} - \mathbf{RES}$  complexity of the input, but there are some that do not correspond to the  $\mathbf{NG} - \mathbf{RES}$ . For that reason, we focus in this paper on the  $\mathbf{C} - \mathbf{RES}$  complexity, as did Mitchell in [32].

Given an instance  $\mathcal{I}$  of a CSP in which every variable has domain  $\{1,...,d\}$ , we construct a boolean CNF-formula CNF( $\mathcal{I}$ ) as follows. For each variable x of  $\mathcal{I}$ , there are d variables in CNF( $\mathcal{I}$ ), denoted x:1,x:2,...,x:d, and there is a domain clause  $(x:1\vee...\vee x:d)$ . For each restriction  $(\delta_1,...,\delta_k)$  on variables  $(x_1,...,x_k)$  in any constraint of  $\mathcal{I}$ , CNF( $\mathcal{I}$ ) has a conflict clause  $(\overline{x_1}:\delta_1\vee...\vee\overline{x_k}:\delta_k)$ . It is easy to see that CNF( $\mathcal{I}$ ) has a satisfying assignment iff  $\mathcal{I}$  does - if  $\mathcal{I}$  has a satisfying assignment, then we produce one for CNF( $\mathcal{I}$ ) by setting  $x:\delta$  to True iff  $x=\delta$ ; if CNF( $\mathcal{I}$ ) has a satisfying assignment, then we produce one for  $\mathcal{I}$  by setting  $x=\delta$  where  $\delta$  is any one of the values for which  $x:\delta$  is True.

**Remark:** It is natural to consider adding an extra set of constraints for each variable x which specify that  $x : \delta$  can be true for at most one value of  $\delta$ . But it is easily verified that each of the results in this paper (in particular, Lemma 7) holds regardless of whether we include these clauses; to be specific, we do not include them.

We define the resolution complexity of  $\mathcal{I}$ , denoted  $\mathbf{C} - \mathbf{RES}(\mathcal{I})$  to be equal to  $\mathbf{RES}(\mathrm{CNF}(\mathcal{I}))$ . In most previous papers bounding the resolution complexity of random instances of SAT or CSP for  $k \geq 3$ , a key lemma has been to establish that the following two conditions hold almost surely for some constants  $\alpha, \zeta > 0$ :

- (A) Every subproblem on at most  $\alpha n$  variables is satisfiable.
- (B) Every subproblem on v variables where  $\frac{1}{2}\alpha n \leq v \leq \alpha n$  has at least  $\zeta n$  variables of degree 1.

For SAT, these two facts imply that a.s. the resolution complexity is exponential in n using principles introduced in [12] and refined to easily applied tools in [8] and [9]. For more general instances of CSP, one needs to establish an additional fact (which is trivially true for SAT):

(C) If x is a variable of degree 1 in a CSP f then, letting f' be the subproblem obtained by removing x and its constraint, any satisfying assignment of f' can be extended to a satisfying assignment of f by assigning some value to x.

In our setting, (C) holds if t < d, but for  $t \ge d$ , it is easy to see that it fails: Suppose that the constraint x lies in contains the d restrictions: (1, 1, 1, ..., 1, 1), (1, 1, 1, ..., 1, 2), ..., (1, 1, 1, ..., 1, d), where x is the last variable in the constraint. Then any satisfying assignment for f' in which all the other variables of the clause receive 1 cannot be extended to f. For  $k \ge 3$  Mitchell's proof[32] applies precisely to the range of f for which (C) holds. For f in which lies the conditions, replacing "degree 1" by "degree 2" in (B) and (C); this revised condition (C) holds precisely when f in the condition f

For higher values of t, we need to replace "degree 1" in condition (B) by a more complicated notion, and then prove that something similar to condition (C) still holds. We describe how to do this in the next section, after presenting several necessary definitions.

#### 2.3 Some new boundaries

A constraint C on variables  $x_1, ..., x_k$  forbids  $x_i : \delta$  if each of the  $d^{k-1}$  possible k-tuples  $(\delta_1, ..., \delta_k)$  with  $\delta_i = \delta$  is a restriction of C. Such a C is called a forbidder. As explained in [5] (and expanded on in [34]), it is the presence of forbidders that causes  $\mathbf{C} - \mathbf{RES}(\mathcal{F}_{n,M=cn}^{d,k,t}) = O(1)$  a.s. for all c when  $t \geq d^{k-1}$ . C permits  $(x_i : \delta, x_j : \gamma)$  if at least one of the  $d^{k-2}$  possible k-tuples  $(\delta_1, ..., \delta_k)$  with  $\delta_i = \delta$  and  $\delta_j = \gamma$  is not a restriction of C. C is a  $(x_i : \delta) \to (x_j : \gamma)$  forcer if C does not permit  $(x_i : \delta, x_j : \gamma')$  for any  $\gamma' \neq \gamma$ ; i.e. if each of the  $(d-1)d^{k-2}$  possible k-tuples  $(\delta_1, ..., \delta_k)$  with  $\delta_i = \delta$  and  $\delta_j \neq \gamma$  is a restriction of C. Thus C implies "If  $x_i = \delta$  then  $x_j = \gamma$ ." In this case, we say that the forcer C starts at  $x_i$  (or, more specifically,  $x_i : \delta$ ) and finishes at  $x_j$  (or  $x_j : \gamma$ ). As predicted by Mitchell[33], it is the presence of forcers that causes  $\mathbf{C} - \mathbf{RES}(\mathcal{F}_{n,M=cn}^{d,k,t}) = \text{poly}(n)$  a.s. for large c when  $(d-1)d^{k-2} \leq t < d^{k-1}$ .

A path of length r in a k-uniform hypergraph H is a sequence of r edges  $e_1, e_2, ..., e_r$  such that:

- for  $1 \le i \le r 1$ ,  $e_i \cap e_{i+1} = x_i$  this is called a *connecting vertex*;
- for all 1 < i < r 2 and j > i + 1,  $e_i \cap e_j = \emptyset$ .
- there are specified vertices  $x_0 \in e_1$  and  $x_r \in e_r$ , called the *endpoints* of the path.

If  $e_1, \ldots, e_r$  is a path and there is an edge  $e_0 \in H$  whose intersection with the vertices of  $e_1, \ldots, e_r$  is only  $\{x_0, x_r\}$  then  $e_0, \ldots, e_r$  form a cycle in H.

A pendant path is a path in which no vertices other than the endpoints lie in any edges of H off the path. In other words, there is no restriction on the degrees of the endpoints, each connecting vertex has degree 2 in H, and every other vertex in the path has degree 1 in H.

A (pendant) path of length r in a CSP is a sequence of r constraints whose underlying edges form a (pendant) path of length r in the underlying hypergraph. If there are values  $\delta_0, ..., \delta_r$  such that the constraint on each  $e_i$  is a  $(x_{i-1}:\delta_{i-1}) \to (x_i:\delta_i)$  forcer, then we say that the (pendant) path is a  $(x_0:\delta_0) \to (x_r:\delta_r)$  forcing (pendant) path. It will be convenient to consider a single variable x to be a forcing path of length zero (note that we trivially have  $(x:\delta) \to (x:\delta)$  for every  $\delta$ ); in this case, both endpoints of the forcing path are considered to be x.

A constraint on the edge  $e_i$  of a path is a P-forcer if it is a  $(x_{i-1}:\delta) \to (x_i,\gamma)$  forcer or a  $(x_i:\gamma) \to (x_{i-1},\delta)$  forcer for some  $\delta,\gamma$ . For any  $\mathcal{I} \in \Omega^{d,k,t}$ :

- The first boundary of  $\mathcal{I}$ , denoted by  $\mathcal{B}^1(\mathcal{I})$ , is the set of non-forbidding constraints of  $\mathcal{I}$  which contain at most one variable of degree greater than 1.
- The second boundary of  $\mathcal{I}$ , denoted by  $\mathcal{B}^2(\mathcal{I})$ , is the set of pendant paths of length 4 in  $\mathcal{I}$  which have no P-forcers.
- The third boundary of  $\mathcal{I}$ , denoted by  $\mathcal{B}^3(\mathcal{I})$ , is the set of pendant paths of length 2 where one of the two constraints is a P-forcer that starts at the connecting vertex, and the other is not a P-forcer that finishes at the connecting vertex.

The boundary of  $\mathcal{I}$  is  $\mathcal{B}(\mathcal{I}) = \mathcal{B}^1(\mathcal{I}) \cup \mathcal{B}^2(\mathcal{I}) \cup \mathcal{B}^3(\mathcal{I})$ . Our main lemma corresponds to conditions (A) and (B) from Section 2.2:

**Lemma 7** Consider any  $\mathcal{I} \in \Omega^{d,k,t}$  on n variables, where  $t < d^{k-1}$ . If for some  $\alpha, \zeta > 0$ , we have:

- (a) every subproblem on at most αn variables is satisfiable, and
- (b) every subproblem  $\mathcal{I}'$  on v variables where  $\frac{1}{2}\alpha n \leq v \leq \alpha n$  has  $|\mathcal{B}(\mathcal{I}')| \geq \zeta n$ ,

then  $\mathbf{C} - \mathbf{RES}(\mathcal{I}) \geq 2^{\Theta(n)}$ .

To prove Lemma 7, we require the following lemma, which corresponds to condition (C) from Section 2.2.

**Lemma 8** Consider any  $\mathcal{I} \in \Omega^{d,k,t}$ , where  $t < d^{k-1}$  and any  $X \in \mathcal{B}(\mathcal{I})$ . Any satisfying assignment of  $\mathcal{I} - X$  can be extended to a satisfying assignment of  $\mathcal{I}$ .

**Proof:** Suppose  $X \in \mathcal{B}^1(\mathcal{I})$ . Then the lemma follows from the fact that since  $t < d^{k-1}$ , X cannot be a forbidding constraint.

Suppose  $X \in \mathcal{B}^2(\mathcal{I})$ , and consider any satisfying assignment of  $\mathcal{I} - X$  where  $x_0, x_4$  are assigned  $\delta_0, \delta_4$ .  $(\mathcal{I} - X)$  is the subproblem obtained by removing all clauses of X and all variables of X other than the endpoints.) Since  $e_1$  is not an  $(x_0, x_1)$ -forcer, there are at least two choices for  $\delta_1$  such that  $e_1$  permits  $(x_0 : \delta_0, x_1 : \delta_1)$ . Similarly, there are at least two choices for  $\delta_3$  which can be assigned to  $x_3$  so  $e_4$  permits  $(x_3 : \delta_3, x_4 : \delta_4)$ . We will show that for at least one of these four choices for the pair  $(\delta_1, \delta_3)$ , there is a value  $\delta_2$  such that  $e_2$  permits  $(x_1 : \delta_1, x_2 : \delta_2)$  and  $e_3$  permits  $(x_2 : \delta_2, x_1 : \delta_3)$ . If this were not the case, then for every  $\delta_2 \in \{1, ..., d\}$  either (i)  $e_2$  does not permit  $(x_1 : \delta_1, x_2 : \delta_2)$  for either of the two choices of  $\delta_1$  (this requires  $2d^{k-2}$  restrictions) or (ii)  $e_3$  does not permit  $(x_2 : \delta_2, x_3 : \delta_3)$  for either of the two choices of  $\delta_3$  (this also requires  $2d^{k-2}$  restrictions). Thus  $e_2, e_3$  would have a total of at least  $2d^{k-1}$  restrictions which is not possible by hypothesis.

Suppose  $X \in \mathcal{B}^3(H)$ . Let the endpoints of X be  $x_0, x_2$ , the connecting variable of X be  $x_1$  and the constraints of X be  $C_1, C_2$  where  $C_1$  is a forcer starting at  $x_1$  and ending at  $x_0$  and  $C_2$  is not a

forcer that starts at  $x_2$  and ends at  $x_1$ . Consider any satisfying assignment of  $\mathcal{I} - X$  where  $x_0, x_2$  are assigned  $\delta_0, \delta_2$ . There are at least d-1 choices for  $\delta_1$  such that  $C_1$  permits  $(x_0 : \delta_0, x_1 : \delta_1)$  and there are at least 2 choices for  $\delta_1$  such that  $C_2$  permits  $(x_1 : \delta_1, x_2 : \delta_2)$ . At least one choice for  $\delta_1$  lies in the intersection of these sets and so the lemma follows.

With Lemma 8 in hand, the proof of Lemma 7 is straightforward, following Mitchell's framework [32].

**Proof of Lemma 7:** Consider any resolution refutation of  $CNF(\mathcal{I})$ . Mitchell ([32], Lemma 1) proves that hypothesis (a) implies there must be a clause C in the refutation and a subproblem  $\mathcal{J}$  of  $\mathcal{I}$  on between  $\frac{1}{2}\alpha n$  and  $\alpha n$  variables, such that  $\mathcal{J}$  minimally implies C in the following sense: (i) Every satisfying assignment of  $\mathcal{J}$  satisfies C, and (ii) For subproblem  $\mathcal{J}'$  of  $\mathcal{J}$ , there is a satisfying assignment of  $\mathcal{J}'$  that does not satisfy C.

We will next prove that C must have at least  $\zeta n/4$  variables:

Consider any clause  $X \in \mathcal{B}^1(\mathcal{J})$ ; we will show that some variable of X appears in C. To see this, consider any assignment  $\alpha$  which satisfies  $\mathcal{J} - X$  but does not satisfy C. By Lemma 8, it is possible to extend  $\alpha$  to a satisfying assignment  $\alpha'$  of  $\mathcal{J}$ , and since  $\mathcal{J}$  implies C,  $\alpha'$  satisfies C. Thus, there is some variable of C that is assigned a value in  $\alpha'$  but not in  $\alpha$ ; that variable must be in X.

Nearly identical arguments show that C must contain a non-endpoint variable of every member of  $\mathcal{B}^2(\mathcal{J})$ , and that C must contain the connecting variable from every member of  $\mathcal{B}^3(\mathcal{J})$ . No variable can be a non-endpoint variable of more than four members of  $\mathcal{B}^2(\mathcal{J})$ . So C contains at least  $|\mathcal{B}^1(\mathcal{J})| + |\mathcal{B}^2(\mathcal{J})|/4 + |\mathcal{B}^3(\mathcal{J})|$  variables. Since by hypothesis (b),  $|\mathcal{B}(\mathcal{J})| \geq \zeta n$ , C must contain at least  $\zeta n/4$  variables, as required.

This allows us to apply the now standard "width lemma" of Ben-Sasson and Wigderson ([9], Corollary 3.6) to prove our lemma. In particular, if we let  $w_1 = k$  be the maximum clause size in  $CNF(\mathcal{I})$ , and  $w_2 \geq \zeta n/4$  be the minimum over all resolution refutations of  $CNF(\mathcal{I})$  of the maximum clause size in the refutation, then the width lemma states that:

$$\mathbf{C} - \mathbf{RES}(\mathcal{I}) = e^{\Omega((w_2 - w_1)^2/n)} \ge 2^{\Theta(n)}.$$

We close this section with a lemma explaining the significance of  $c^*(d, k, t)$ .

**Lemma 9** For  $d, k \geq 2$  and  $t \geq (d-1)d^{k-2}$ , let  $c = c^*(d, k, t) = \frac{1}{dk(k-1)} {d^k \choose t} / {d^k - (d-1)d^{k-2} \choose t - (d-1)d^{k-2}}$ . Specify any variable x and value  $\delta \in \{1, ..., d\}$ . The expected number of forcers in  $\mathcal{F}_{n, M=cn}^{d, k, t}$  starting with  $x : \delta$  is 1.

**Proof:** Let  $L = (d-1)d^{k-2}$ . The expected number of constraints containing x is ck. For each of the d(k-1) choices of  $x' \neq x, \delta' \in \{1, ..., d\}$  for a particular constraint containing x, the probability that the constraint forms a  $(x:\delta) \to (x':\delta')$  forcer is  $\binom{d^k-L}{t-L}/\binom{d^k}{t}$ , as there are exactly  $\binom{d^k-L}{t-L}$  choices for such a forcer. Therefore, the expected number of forcers from  $x:\delta$  is

$$cdk(k-1) \times \frac{\binom{d^k-L}{t-L}}{\binom{d^k}{t}} = \frac{\binom{d^k}{t}}{dk(k-1)\binom{d^k-L}{t-L}} \times dk(k-1) \times \frac{\binom{d^k-L}{t-L}}{\binom{d^k}{t}} = 1$$

It is instructive to consider the case k=d=2 and  $t=(d-1)d^{k-2}=1$  which is the more familiar random 2-SAT. Here, every 2-clause can be viewed as the union of two forcers, eg.  $(x_1 \vee x_2)$ 

is equivalent to the conjunction of the two forcers  $(x_1:F) \to (x_2:T)$  and  $(x_2:F) \to (x_1:T)$ . (Of course, we are considering the domain to be  $\{T(\text{rue}),F(\text{alse})\}$  rather than  $\{1,2\}$ . Note that  $c^*=1$  which is the satisfiability threshold for 2-SAT. The reader who is familiar with random 2-SAT, will recognize that the property guaranteed by Lemma 9 corresponds very closely to what happens to cause the random 2-SAT to be unsatisfiable. Thus, it is not surprising that for general d, k, at  $c > c^*$  the forcers alone produce an unsatisfiable formula and that it, like random 2-SAT, has small resolution complexity.

For d=2 and general k, it is easy to see that an  $(x:F) \to (y:T)$  forcer is also an  $(y:F) \to (x:T)$  forcer. Thus, such a forcer implies the 2-clause  $(x \vee y)$ . Extending this reasoning shows that for  $c>c^*$ , the forcers alone will contain a random instance of 2-SAT where the number of 2-clauses is above the satisfiability threshold. As mentioned earlier, this was discovered independently by Gao and Culberson [23].

### 3 Proof of Theorem 1

We begin with a lemma of a type that has become standard in papers on the resolution complexity of random formulae. It says that a.s. every subproblem of  $\mathcal{F}_{n,M}^{d,k,t}$  with at most  $\alpha n$  variables has a very low clause-vertex ratio. Thus, to prove that the conditions of Lemma 7 hold, it suffices to prove that certain types of subproblems must have high clause-vertex ratio.

**Lemma 10** Let c > 0 and  $k \ge 2$ , and let H be the random k-uniform hypergraph with n vertices and m = cn edges. Then for any  $\delta > 0$ , there exists  $\alpha = \alpha(c, k, \delta) > 0$  such that a.s. H has no subgraph with  $0 < h \le \lfloor \alpha n \rfloor$  vertices and at least  $\left(\frac{1+\delta}{k-1}\right)h$  edges.

**Proof:** This proof follows a straightforward first moment calculation of a type that has been carried out many times in similar settings, starting with Lucsak[31].

Let  $\mu = \frac{1+\delta}{k-1}$ ,  $S_h$  be the number of subgraphs of H with h vertices and  $exactly \lceil \mu h \rceil$  edges, and set  $S = \sum_{h=1}^{\lfloor \alpha n \rfloor} S_h$ . Note that if  $S_h = 0$  then there are no subgraphs of H with h vertices and at least  $\mu h$  edges. To count  $\mathbf{E}(S_h)$ , we multiply the number of choices of h vertices by the number of choices for  $\lceil \mu h \rceil$  of the cn random edges and the probability that each of those random edges lies entirely in that set of h vertices:

$$\mathbf{E}(\mathcal{S}_h) \le \binom{n}{h} \binom{cn}{\lceil \mu h \rceil} \left(\frac{h}{n}\right)^{k\lceil \mu h \rceil}.$$

This yields:

$$\mathbf{E}(\mathcal{S}) = \sum_{h=1}^{\lfloor \alpha n \rfloor} \mathbf{E}(\mathcal{S}_h)$$

$$\leq \sum_{h=1}^{\lfloor \alpha n \rfloor} \binom{n}{h} \binom{cn}{\lceil \mu h \rceil} \left(\frac{h}{n}\right)^{k\lceil \mu h \rceil}$$

$$\leq \sum_{h=1}^{\lfloor \alpha n \rfloor} \left(\frac{en}{h}\right)^h \left(\frac{ecn}{\lceil \mu h \rceil}\right)^{\lceil \mu h \rceil} \left(\frac{h}{n}\right)^{k\lceil \mu h \rceil}$$

$$\leq \frac{ec}{\mu} \sum_{h=1}^{\lfloor \alpha n \rfloor} \left[\left(\frac{h}{n}\right)^{(k-1)\mu-1} e^{\mu+1} (c/\mu)^{\mu}\right]^h$$

$$= \frac{ec}{\mu} \sum_{h=1}^{\lfloor \alpha n \rfloor} \left[ \left( \frac{h}{n} \right)^{\delta} c' \right]^{h} \qquad \text{for } c' = e^{\mu+1} (c/\mu)^{\mu}$$

$$\leq \frac{ec}{\mu} \left( \sum_{h=1}^{\lfloor \log n \rfloor} \left[ \left( \frac{\lfloor \log n \rfloor}{n} \right)^{\delta} c' \right] + \sum_{h=\lfloor \log n \rfloor+1}^{\lfloor \alpha n \rfloor} \left[ \left( \frac{\lfloor \alpha n \rfloor}{n} \right)^{\delta} c' \right]^{\lfloor \log n \rfloor} \right)$$

$$\leq O\left( \frac{\log^{1+\delta} n}{n^{\delta}} \right) + \sum_{h=\lfloor \log n \rfloor+1}^{\lfloor \alpha n \rfloor} O\left( \frac{1}{n^{2}} \right) \qquad \text{for sufficiently small } \alpha$$

$$= o(1).$$

The next lemma will allow us to show that subproblems with small boundaries must have high clause-variable ratio, and hence, by the previous lemma, must be large.

**Lemma 11** Let  $r \geq 2$  be a constant, and H be a k-uniform hypergraph on n vertices and m edges that does not have any component which is a cycle. Let  $B_1$  be the set of edges which have at most one vertex of degree greater than 1, and  $B_2$  be the set of pendant paths of length r. If  $|B_1| + |B_2| \leq n/(72r^2k^3)$ , then for  $\delta = \frac{1}{3rk^2}$ :  $m \geq n\left(\frac{1+\delta}{k-1}\right)$ .

Intuitively, the lemma is clear: If  $|B_1| = 0$  and if no vertex has degree greater than 2, then it is easy to see that we would have m = n/(k-1). So in order for m to not be much bigger, either  $B_1$  must be large or there must be very few vertices of degree greater than 2. If the latter is true and  $B_1$  is small, then H must contain long pendant paths and so  $B_2$  will be big. Our formal proof is a bit lengthy and so we defer it to the end of the section.

These two lemmas are enough to prove the first of our main theorems:

**Proof of Theorem 1:** It suffices to prove that a.s. conditions (a) and (b) of Lemma 7 hold for  $\mathcal{F}_{n,M=cn}^{d,k,t}$ , where  $\alpha = \alpha(c,k,\delta=1/(12k^2))$  from Lemma 10 and  $\zeta = \min(1/(72\times 16k^3),\alpha/10k)$ .

Since  $t < (d-1)d^{k-2}$ ,  $\mathcal{F}_{n,M}^{d,k,t}$  has no forcer constraints. Using this fact, our proof would follow immediately from Lemmas 10 and 11, if it were not for the fact that Lemma 11 only applies to hypergraphs with no cycle components.

We begin with condition (a). Suppose that  $\mathcal{J}$  is a minimally unsatisfiable subproblem of  $\mathcal{F}_{n,M}^{d,k,t}$ . Thus, the underlying hypergraph of  $\mathcal{J}$  is connected. Furthermore, since  $t < (d-1)d^{k-2}$ , it is easily verified that the underlying hypergraph of  $\mathcal{J}$  cannot be a single cycle. Finally, Lemma 8 implies that  $|\mathcal{B}^1(\mathcal{J})| = |\mathcal{B}^2(\mathcal{J})| = 0$ . Therefore, since  $\mathcal{J}$  has no forcers, Lemma 11 with r = 4 applies to the underlying hypergraph of  $\mathcal{J}$  and so  $\mathcal{J}$  has clause-variable ratio at least  $(1+\delta)/(k+1)$ . Thus Lemma 10 implies that a.s.  $\mathcal{F}_{n,M}^{d,k,t}$  has no minimally unsatisfiable subproblems of size at most  $\alpha n$ . Therefore a.s.  $\mathcal{F}_{n,M}^{d,k,t}$  has no unsatisfiable subproblems of size at most  $\alpha n$ .

Therefore a.s.  $\mathcal{F}_{n,M}^{d,k,t}$  has no unsatisfiable subproblems of size at most  $\alpha n$ . Next, condition (b). We will use the easy fact, provided as Lemma 22 in Appendix B, that a.s. the underlying random hypergraph of  $\mathcal{F}_{n,M}^{d,k,t}$  has fewer than  $\log n$  cycles of length at most 4. Suppose, by contradiction, that  $\mathcal{J}$  is a subproblem of  $\mathcal{F}_{n,M}^{d,k,t}$  with v variables where  $\frac{1}{2}\alpha n \leq v \leq \alpha n$ , and with  $|\mathcal{B}^1(\mathcal{J})| + |\mathcal{B}^2(\mathcal{J})| \leq \zeta n$ . (Since there are no forcers,  $|B^3(\mathcal{J})| = 0$ .) Let H' be the subhypergraph obtained by removing all the cycle components from the underlying hypergraph of  $\mathcal{J}$ . By Lemma 11 (with r=4), H' has at least  $|H'|\frac{1+\delta}{k-1}$  edges, and note also that  $|H'| \leq |\mathcal{J}| \leq \zeta n < \alpha n$ . By Lemma 10, a.s. every such H' is empty. Thus a.s. for every such subproblem  $\mathcal{J}$ , every component in the underlying hypergraph of  $\mathcal{J}$  is a cycle. Every vertex in such cycle of length at least 5 must lie in a member of  $\mathcal{B}^2(\mathcal{J})$ , and every member of  $\mathcal{B}^2(\mathcal{J})$  contains fewer than 4k vertices; so there are at most  $4k\zeta n$  vertices in those cycles. As mentioned above, a.s. there are at most  $4\log n$  vertices which lie in cycles of length at most 4 in  $\mathcal{F}_{n,M}^{d,k,t}$ . Since  $4k\zeta n + 4\log n < \frac{1}{2}\alpha n \leq |\mathcal{J}|$ , we have a contradiction.

We now close this section with the proof of Lemma 11.

**Proof of Lemma 11:** We say that a pendant path p of length at least 1 is *contractible* if (i) both of it's endpoints have degree 2 and (ii) p is maximal in the sense that it is not part of a longer pendant path whose endpoints both have degree 2. Let  $\mathcal{P}$  be the set of contractible pendant paths in H.

We form a hypergraph H' with no long pendant paths as follows:

For each  $p \in \mathcal{P}$  with endpoints x, y, we remove all edges and vertices of p except for x, y and do the following: If x, y do not both lie in some edge outside of p, then we contract x, y into a single vertex. Otherwise we create a new edge containing x, y and k-2 new degree 1 vertices; this edge is called a reduced edge and x, y are its endpoints. Note that, since H has no cycle components, this operation does not create any new contractible paths, and it does not destroy any contractible paths. So there is no need to iterate this process, and we contract every path from  $\mathcal{P}$ . For convenience, we first contract all paths of length at least r in Phase A, and then all the remaining paths in Phase B.

The resulting hypergraph is H'. Every non-reduced edge in H' either (i) has a vertex of degree at least 3, (ii) has at least 3 degree 2 vertices or (iii) is in  $B_1$ .

n' is the number of vertices in H'; s is the number of vertices of degree greater than 1; m' is the number of edges in H';  $m_2$  is the number of reduced edges;  $m^*$  is the number of edges with at least one degree 3 vertex.

Since H has no cycle components, the endpoints of any reduced edge must lie in a non-reduced edge. Since those endpoints have degree 2, each edge can hold the endpoints of at most k/2 reduced edges. Therefore,  $m_2 \leq m' \cdot \frac{k}{2}/(\frac{k}{2}+1) = m'(k/(k+2))$ .

The sum of the degrees of the vertices is n'+s plus the sum of all vertices v with  $\deg(v) \geq 3$  of  $\deg(v)-2$ . Each such vertex v lies in  $\deg(v)$  edges, and so by counting  $(\deg(v)-2)/\deg(v) \geq \frac{1}{3}$  for each of those edges, it follows that this latter sum is at least  $m^*/3$ . Therefore, since the sum of the degrees of all vertices is km', we have:

$$km' \ge n' + s + m^*/3.$$

Furthermore, by counting the number of degree 1 variables in each edge, we have:

$$s \ge n' - (k-3)m' - (m^* + m_2) - 2|B_1|.$$

These two equations combine to yield:

$$km' \geq 2n' - (k-2)m' + (m' - 2m^*/3 - m_2) - 2|B_1|$$
  
 
$$\geq 2n' - (k-2)m' + \frac{1}{3}(m' - m_2) - 2|B_1|$$
  
 
$$\geq 2n' - m'\left(k - 2 - \frac{2}{3(k+2)}\right) - 2|B_1|,$$

and so  $m' \ge (2n'-2|B_1|)/(2k-2-\frac{2}{3(k+2)}) \ge n'\left(\frac{1+2\delta}{k-1}\right)-2|B_1|$ , since  $\delta < \frac{1}{3k(k+2)}$  and  $2k-2-\frac{2}{3(k+2)} > 1$ .

Now we observe that very few vertices were removed during Phase A. Any pendant path of length  $l \ge r$  contains at least l - r + 1 pendant paths of length r. So the total of the lengths of all

such paths in H is at most  $r|B_2|$ . Therefore, the number of vertices in the hypergraph at the end of Phase A is at least  $n - rk|B_2|$ .

Now we consider what happened during Phase B. Every time we contracted a contractible path p of length l < r to a vertex, the net loss in edges was l and the net loss of vertices was (k-1)l-1, and each time we contracted one to an edge, those net losses were l-1 and (k-1)(l-1)-1. Thus, for some v, we lost v vertices and at least  $v(r-1)/((r-1)(k-1)-1) > v(1+\frac{1}{rk})/(k-1) \ge v\left(\frac{1+2\delta}{k-1}\right)$  edges during Phase B, since  $\delta \le \frac{1}{2rk}$ .

Therefore, the hypergraph remaining at the end of Phase A has at least

$$(n - rk|B_2|) \left(\frac{1 + 2\delta}{k - 1}\right) - 2|B_1| \ge n\left(\frac{1 + \delta}{k - 1}\right) + \frac{n\delta}{k - 1} - 3r|B_2| - 2|B_1| > n\left(\frac{1 + \delta}{k - 1}\right)$$

edges. Since Phase A does not increase the number of edges, this proves our Lemma.

#### 4 Proof of Theorem 2

A  $Z_q$  configuration is a collection of q vertex-disjoint forcing paths in  $\mathcal{I}$ , each with possibly length zero (i.e. a single vertex), plus  $q\left(\frac{1+\gamma}{k-1}\right)$  other edges, each containing k endpoints of the paths, where  $\gamma = 1/(300k^2)$ .

For  $\mathcal{I} \in \Omega^{d,k,t}$  and  $t \geq (d-1)k^{d-2}$ , let  $\mathcal{P} = p_1, \ldots, p_q$  be a collection of forcing paths of  $\mathcal{I}$  such that: (i) every vertex lies in exactly one of these paths, (ii) it is not possible to transform  $\mathcal{P}$  into a collection of q-1 paths meeting condition (i) by adding another forcer from  $\mathcal{I}$ . Obviously a collection of paths of length 0, one for each variable of  $\mathcal{I}$ , satisfies (i), and so some collection exists which satisfies (i) and (ii).

**Lemma 12** For any  $\mathcal{I}$ ,  $\mathcal{P}$  as described above: If  $|\mathcal{B}(\mathcal{I})| < q/(72000k^3)$  and if the underlying hypergraph of  $\mathcal{I}$  has no cycle components then  $\mathcal{I}$  contains a  $Z_q$  configuration.

**Proof:** Suppose that among the paths in  $\mathcal{P}$ , exactly  $p_1, \ldots, p_r$   $(r \leq q)$  have at least one edge each; the others have length 0. Consider a path  $p_i$ ,  $(1 \leq i \leq r)$ , and suppose it has  $l \geq 1$  edges. Let  $x_0$  and  $x_l$  be the start and end points of this path; so  $p_i$  is a  $(x_0:\delta_0) \to (x_l:\delta_l)$  forcer, for some values  $\delta_0, \delta_l$ . Remove from  $\mathcal{I}$  all of the clauses and variables of  $p_i$  other than  $x_0, x_l$ . Then add k-2 new variables and a  $(x_0:\delta_0) \to (x_l:\delta_l)$  forcer. The new forcer is called a reduced forcer. We do this for every  $p_i, 1 \leq i \leq r$ . The new CSP obtained after these operations is denoted by  $\mathcal{I}'$ . Note that n', the number of variables in  $\mathcal{I}'$ , is q + r(k-1). Note also that since  $p_1, \ldots, p_q$  are vertex-disjoint, no two reduced forcers in  $\mathcal{I}'$  share a vertex.

Claim 13 There is no forcer path of length at least 4 in  $\mathcal{I}'$ .

**Proof:** By contradiction, assume that  $p = e_1 e_2 e_3 e_4$  is a forcer path in  $\mathcal{I}'$ , with  $u_{i-1}, u_i$  being the start and end points of  $e_i$ . If  $e_2$  is not a reduced forcer, then adding the forcer  $e_2$  to  $\mathcal{P}$  would concatenate the path in  $\mathcal{P}$  containing  $u_1$  with the one containing  $u_2$  without violating condition (i). This contradicts condition (ii). If  $e_2$  is a reduced forcer, then  $e_3$  cannot be (since no two reduced forcers share a variable). Thus a similar contradiction arises when we consider adding the forcer  $e_3$  to  $\mathcal{P}$ .

**Claim 14** Every pendant path of length 10 in  $\mathcal{I}'$  has a subpath which is in  $\mathcal{B}^2(\mathcal{I}) \cup \mathcal{B}^3(\mathcal{I})$ .

**Proof:** Let P' be a pendant path in  $\mathcal{I}'$ . By replacing each reduced edge of P' by its corresponding path from  $\mathcal{P}$ , we obtain a path P in  $\mathcal{I}$ . Every non-forcer in P' is a non-forcer in P. Assume  $P = e_1, \ldots, e_l$  where  $x_i = e_i \cap e_{i+1}$ , and that  $e_a$  and  $e_b$ , b > a, are two non-forcers in P', and hence in P, such that there is no other non-forcer between them. If b - a > 1 then  $e_{a+1}$  is a forcer, and it must start at  $x_{a+1}$  for otherwise  $\{e_a, e_{a+1}\} \in \mathcal{B}^3(\mathcal{I})$  and we are done. Similarly,  $e_{b-1}$  must be a forcer starting at  $x_{b-2}$ , otherwise  $\{e_{b-1}, e_b\} \in \mathcal{B}^3(\mathcal{I})$ . But these two imply that along the path  $e_{a+1}, \ldots, e_{b-1}$  there is a member of  $\mathcal{B}^3(\mathcal{I})$ . Thus, we can assume that b - a = 1, i.e. that the non-forcers in P' are consecutive. Also, if i is the largest index for which  $e_i$  is a non-forcer in P, then  $e_l e_{l-1} \ldots e_{i+1}$  must be a forcing path going into  $x_i$ , for otherwise there is a member of  $\mathcal{B}^3(\mathcal{I})$  along this path. Therefore the portions of P' on the sides of these non-forcers form forcing paths. Similar arguments show that if P' has no non-forcers, then it contains at most two forcing paths starting at the end points of P', or else P' contains a member of  $\mathcal{B}^3(\mathcal{I})$ .

By Claim 13, the length of each forcing path is at most 3. Therefore, if P' has length 10, then P' has at least 4 consecutive non-forcers and so that subpath of P' is in  $\mathcal{B}^2(\mathcal{I})$ .

Let  $B_1$  be the set of clauses which have at most 1 variable of degree greater than 1 in  $\mathcal{I}'$  and  $B_2$  be the set of pendant paths of length 10 in  $\mathcal{I}'$ . Note that  $|\mathcal{B}^1(\mathcal{I})| \geq |B_1|$ . Since no subpath can lie in more than 10 members of  $B_2$ , Claim 14, implies that  $|B_2| \leq 10(|\mathcal{B}^2(\mathcal{I}) + \mathcal{B}^3(\mathcal{I})|)$ . Therefore,  $|B_1| + |B_2| \leq 10|\mathcal{B}(\mathcal{I})| \leq q/(7200k^3) < n'/(7200k^3)$ , as n', the number of variables in  $\mathcal{I}'$ , is q + r(k-1). So applying Lemma 11 with r = 10 to the underlying hypergraph of  $\mathcal{I}'$ , the number of clauses in  $\mathcal{I}'$  is at least  $n'\left(\frac{1+\gamma}{k-1}\right)$ , and at least  $n'\left(\frac{1+\gamma}{k-1}\right) - r \geq q\left(\frac{1+\gamma}{k-1}\right)$  are not reduced forcers. Those clauses and the paths in  $\mathcal{P}$  will form a  $Z_q$  configuration in  $\mathcal{I}$ .

**Lemma 15** For any constants  $d, k \geq 2$  and  $(d-1)d^{k-2} \leq t < d^{k-1}$ , and for every  $c < c^*(d, k, t)$ , there exists  $\sigma > 0$ , such that a.s.  $\mathcal{F}_{n, M = cn}^{d, k, t}$  has no  $Z_q$  configuration with  $q \leq \sigma n$ .

**Proof:** For this proof, it will be convenient to work in the  $\mathcal{F}_{n,p}^{d,k,t}$  model described in Section 2.1 where each of the  $\binom{n}{k}$  potential edges is chosen for the constraint hypergraph with probability  $ck!/n^{k-1}$ . Lemma 23 in Appendix B shows that proving this model a.s. has no  $Z_q$  configuration with  $q \leq \sigma n$  will imply that a.s. neither does  $\mathcal{F}_{n,M=cn}^{d,k,t}$ . We compute the expected number of  $Z_q$  configurations. To do so, we suppose that the q forcing

We compute the expected number of  $Z_q$  configurations. To do so, we suppose that the q forcing paths are ordered  $p_1, ..., p_q$ , when we count the number of ways to choose them. Since the forcing paths of a  $Z_q$  configuration are actually unordered, this produces an overcount which we correct by dividing by q!.

We start with the computations related to the forcing paths: For each i, let  $a_i \geq 0$  be the number of forcers in  $p_i$ , and set  $A = \sum_{i=1}^q a_i$ . The number of ways to choose the end-points of the paths is at most  $n^{2q}$  (fewer if some of the path lengths are 0). For each  $p_i$ , there are at most  $n^{a_i-1}$  choices for the connecting variables, and  $d^{a_i+1}$  choices of values to use on the variables to form a forcing path. Also, for each edge of each  $p_i$ , there are there are  $\binom{n}{k-2}$  choices for the set of degree 1 variables. The probability that all the specified hyperedges exist is  $(ck!/n^{k-1})^A$ . The probability that the constraints all form the specified forcers is, as we argued in the proof of Lemma 9,  $\binom{d^k-(d-1)d^{k-2}}{t-(d-1)d^{k-2}}/\binom{d^k}{t}$ .

Now we turn our attention to the additional edges: There are at most  $\binom{2q}{k}$  potential edges containing only endpoints of the paths, and we must choose  $B = q\left(\frac{1+\gamma}{k-1}\right)$  of them; there are  $\binom{\binom{2q}{k}}{B}$  ways to do so. The probability that the B chosen edges are all present in H is  $(ck!/n^{k-1})^B$ .

Letting  $c = (1 - \epsilon)c^*$ , this yields:

$$\begin{split} \mathbf{E}(|Z_{q}|) & \leq \frac{n^{2q}}{q!} \times \sum_{a_{1},...,a_{q} \geq 0} n^{A-q} d^{A+q} \binom{n}{k-2}^{A} \left[ \frac{(1-\epsilon)c^{*}k!}{n^{k-1}} \right]^{A} \left[ \frac{\binom{d^{k}-(d-1)d^{k-2}}{t-(d-1)d^{k-2}}}{\binom{d^{k}}{t}} \right]^{A} \\ & \times \binom{\binom{2q}{k}}{B} \binom{ck!}{n^{k-1}}^{B} \\ & \leq \frac{n^{2q}}{q!} \binom{\binom{2q}{k}}{B} \binom{ck!}{n^{k-1}}^{B} \times \sum_{a_{1},...,a_{q} \geq 0} \binom{n}{k-2}^{A} n^{A-q} d^{A+q} \left[ \frac{(1-\epsilon)(k-2)!}{dn^{k-1}} \right]^{A} \\ & \leq \frac{n^{2q}}{q!} \left( \frac{(2q)^{k}e}{k!B} \frac{ck!}{n^{k-1}} \right)^{B} \times \sum_{a_{1},...,a_{q} \geq 0} \frac{n^{(k-2)A+A-q} d^{A+q}}{(k-2)!A} \left[ \frac{(1-\epsilon)(k-2)!}{dn^{k-1}} \right]^{A} \\ & \leq \frac{n^{2q}}{q!} \left( \frac{(2q)^{k}ec}{n^{k-1}(q/k-1)} \right)^{q(1+\gamma)/(k-1)} \times \frac{d^{q}}{n^{q}} \sum_{a_{1},...,a_{q} \geq 0} (1-\epsilon)^{A} \\ & \leq \binom{\psi^{q}}{n}^{\gamma q} \left[ \sum_{i \geq 0} (1-\epsilon)^{i} \right]^{q} \\ & \leq \binom{\psi^{\prime}q}{n}^{\gamma q}, \end{split}$$

where  $\psi, \psi' > 0$  are functions of  $c, d, k, \epsilon, \gamma$ , and hence of c, d, k, t since  $\epsilon, \gamma$  can be derived from c, d, k, t. Therefore,

$$\sum_{q=1}^{\sigma n} \mathbf{E}(|Z_q|) \leq \sum_{q=1}^{\lfloor \log n \rfloor} \left(\frac{\psi'q}{n}\right)^{\gamma q} + \sum_{q=\lfloor \log n \rfloor+1}^{\sigma n} \left(\frac{\psi'q}{n}\right)^{\gamma q}$$

$$\leq O\left(\frac{\log^2 n}{n}\right) + \sum_{q=\lfloor \log n \rfloor+1}^{\sigma n} O(n^{-2}) = o(1),$$

for  $\sigma > 0$  sufficiently small that  $(\psi'\sigma)^{\gamma} < e^{-3}$ . Thus  $\sigma$  is a function of  $c, d, k, t, \gamma$  and hence of c, d, k, t.

**Lemma 16** For  $c < c^*(d, k, t)$ , and for every constant integer  $\theta > 0$ , a.s. the number of variables of  $\mathcal{F}_{n,M=cn}^{d,k,t}$  that are part of a maximal forcing path of length greater than  $\theta$  is at most  $3\theta d(c/c^*)^{\theta}$ .

**Proof:** Assume that  $c = (1 - \epsilon)c^*$ . We will again work in the  $H_{n,p}$  model where  $p = c \times k!/n^{k-1}$ . Lemma 23 from Appendix B permits us to do so.

For any value of  $\theta$ , let  $X_{\theta}$  be the number of forcing paths of length  $\theta$  and let  $Y_{\theta}$  be the number of maximal forcing paths of length at least  $\theta$ . Obviously,  $X_{\theta}$  is an upper bound for  $Y_{\theta}$  (since for every forcing path of length  $\theta$  there is a unique maximal forcing path of length at least  $\theta$ ). Our goal is to upper bound  $Y_{\theta}$  for constant values of  $\theta$  and for that we upper bound  $X_{\theta}$ . We first compute the expected value of  $X_{\theta}$ . We have to choose  $\theta + 1$  variables for the connecting points and the end-points of the path; there are  $\binom{n}{\theta+1}$  ways to do so; We order them; there are  $(\theta+1)!$  ways to do so. Then we choose one of d values for each; there are  $d^{\theta+1}$  ways to do so. Then we choose the remaining k-2 vertices for each of the  $\theta$  constraints; there are  $\binom{n-\theta-1}{(k-2)\theta}\frac{[(k-2)\theta]!}{(k-2)!^{\theta}}$  ways to do so.

Finally, we multiply by the probability that the edge for each constraint is chosen in the underlying hypergraph, and that the random constraints chosen for those edges are the specified forcers; the first of these probabilities is  $p^{\theta}$  and the second is  $(dk(k-1)/c^*)^{\theta}$  (the latter computation uses the same arguments found in the proof of Lemma 9). This yields:

$$\mathbf{E}(X_{\theta}) = \binom{n}{\theta+1} (\theta+1)! d^{\theta+1} \binom{n-\theta-1}{(k-2)\theta} \frac{[(k-2)\theta]!}{(k-2)!^{\theta}} \left(\frac{(1-\epsilon)(k-2)!}{dn^{k-1}}\right)^{\theta}$$

$$= \frac{n!}{[n-\theta-1-(k-2)\theta]!} \cdot \frac{(1-\epsilon)^{\theta}d}{n^{(k-1)\theta}}$$

$$= (1+o(1))nd(1-\epsilon)^{\theta} \quad \text{for constant values of } \theta$$

Now we bound the probability of  $X_{\theta} > 2\mathbf{E}(X_{\theta})$  using the second moment method. By Chebychev's Inequality, this probability is at most  $(\mathbf{E}(X_{\theta}^2) - \mathbf{E}(X_{\theta})^2)/\mathbf{E}(X_{\theta})^2$ . So to prove that this probability is o(1), it will suffice to prove that  $\mathbf{E}(X_{\theta}^2) \leq \mathbf{E}(X_{\theta})^2(1 + o(1))$ .

Consider a fixed forcing path A of length  $\theta$ . For each  $i, j \geq 1$  we bound the number of potential forcing paths B of length  $\theta$  which have exactly i constraints in common with A and for which these i constraints form j segments in A. To compute this, we first choose the j segments by choosing the vertices of their end-points from A and the positions of those endpoints in B; there are  $\binom{\theta}{2j}^2$  ways to make this selection. Then we match the j segments of A to the j segments of B; there are j! choices for this. Then we select  $\theta+1-i-j$  points (the rest of the connecting variables of path B) and a value for each; there are at most  $(nd)^{\theta+1-i-j}$  ways to do this. Then for each of the  $\theta-i$  constraints in B-A we select the remaining k-2 variables; there are  $\binom{n}{(k-2)(\theta-i)}\frac{[(k-2)(\theta-i)!]}{(k-2)!(\theta-i)}$  possible choices. Finally we multiply by the probability that the  $\theta-i$  constraints in B-A are selected and are forcers; as in the previous calculation, this probability is  $(\frac{(1-\epsilon)(k-2)!}{dn^{k-1}})^{\theta-i}$ . For any potential forcing path A of length  $\theta$  let  $Q_A$  be its indicator variable.

$$\begin{split} \mathbf{E}(X_{\theta}^{2}) &= \sum_{A,B} \mathbf{Pr}(Q_{A} = 1 \cap Q_{B} = 1) \\ &= \sum_{A} \sum_{B:A \cap B = \emptyset} \mathbf{Pr}(Q_{A} = 1) \mathbf{Pr}(Q_{B} = 1) + \sum_{A} \sum_{B:A \cap B \neq \emptyset} \mathbf{Pr}(Q_{A} = 1 \cap Q_{B} = 1) \\ &\leq \mathbf{E}(X_{\theta})^{2} + \sum_{A} \mathbf{Pr}(Q_{A} = 1) \times \sum_{i=1}^{\theta} \sum_{j=1}^{i} \left(\frac{\theta}{2j}\right)^{2} j! n^{\theta + 1 - i - j} d^{\theta + 1 - i - j} \\ &\times \left(\binom{n}{(k-2)(\theta-i)}\right) \frac{\left[(k-2)(\theta-i)!\right]}{(k-2)!(\theta-i)} \left(\frac{(1-\epsilon)(k-2)!}{dn^{k-1}}\right)^{\theta-i} \\ &\leq \mathbf{E}(X_{\theta})^{2} + \mathbf{E}(X_{\theta}) \sum_{i=1}^{\theta} \sum_{j=1}^{i} \left(\frac{e^{2}\theta^{2}}{4j^{2}}\right)^{2j} j^{j} n^{1-j} d^{1-j} (1-\epsilon)^{\theta-i} \\ &\leq \mathbf{E}(X_{\theta})^{2} + \mathbf{E}(X_{\theta}) \cdot n d(1-\epsilon)^{\theta} \sum_{i=1}^{\theta} \sum_{j=1}^{i} \left(\frac{\alpha \theta^{4}}{j^{3} n d}\right)^{j} (1-\epsilon)^{-i} \quad \text{(for a constant } \alpha > 0) \\ &\leq \mathbf{E}(X_{\theta})^{2} + \mathbf{E}(X_{\theta})^{2} \cdot (1+o(1)) \cdot O\left(\frac{1}{n}\right). \end{split}$$

Therefore, a.s.  $X_{\theta} \leq 2(1+o(1))nd(1-\epsilon)^{\theta}$ . Thus the number of variables which lie on a maximal forcing path of length at least  $\theta$  is a.s. at most  $2(1+o(1))\theta nd(1-\epsilon)^{\theta} < 3\theta d(c/c^*)^{\theta}$ .

**Proof of Theorem 2:** The proof is nearly identical to that of Theorem 1, this time using Lemmas 12 and 15 rather than Lemmas 11 and 10. We will prove that conditions (a,b) of Lemma 7 hold with  $\alpha = \frac{1}{2}\sigma$  and  $\zeta = \alpha/(4\theta \times 72000k^3)$  where  $\theta$  is a positive integer such that  $3\theta d(c/c^*)^{\theta} \leq \alpha/8$ .

We start with condition (a). Let  $\mathcal{J}$  be a minimally unsatisfiable subproblem of  $\mathcal{F}_{n,M}^{d,k,t}$  and Hbe the underlying hypergraph of  $\mathcal{J}$ . Clearly H is connected. Next we show that H cannot be a single cycle. By way of contradiction suppose that H is a cycle with constraints  $C_1, \ldots, C_\ell$ . Since  $t < d^{k-1}$  we have: for any constraint  $C_i$  and any pair of variable/value  $(x:\delta)$  with  $x \in C_i$  there is a satisfying assignment of values to the variables of  $C_i$  in which x gets value  $\delta$ . Let  $x_1 = C_1 \cap C_2$ ,  $x_{\ell} = C_1 \cap C_{\ell}$ , and consider  $\mathcal{J} - C_1$ . Consider any value  $\delta \in \mathcal{D}$ . There is a satisfying assignment for  $C_1$  in which x gets  $\delta$  - let  $\delta_2$  be the value assigned to  $x_2$  in that assignment. There is a satisfying assignment for  $C_2$  in which  $x_2$  gets  $\delta_2$  - let  $\delta_3$  be the value assigned to  $x_3$  in that assignment. Repeating this argument yields a satisfying assignment for  $\mathcal{J}-C_1$  in which  $x_1$  gets  $\delta$ . Since this is true of every value  $\delta$ , there are at least d pairs  $(\delta, \delta')$  in which there is a satisfying assignment for  $\mathcal{J}-C_1$  in which  $x_1$  gets  $\delta$  and  $x_\ell$  gets  $\delta'$ . Since  $C_1$  contains fewer than  $d\times d^{k-2}$  restrictions, at least one of these pairs  $(\delta, \delta')$  is such that there is a satisfying assignment to  $C_1$  in which  $x_1$  gets  $\delta$  and  $x_{\ell}$  gets  $\delta'$ . Therefore,  $\mathcal{J}$  is satisfiable, which is a contradiction. Thus H is connected and is not a cycle. Since  $\mathcal{J}$  is minimally unsatisfiable, Lemma 8 implies that  $|\mathcal{B}(\mathcal{J})| = 0$ . So by Lemma 12,  $\mathcal{J}$  has a  $Z_q$  configuration and by Lemma 15, that configuration has  $q \geq \sigma n$ . Since the paths of the  $Z_q$  configuration are vertex-disjoint, we have  $|\mathcal{J}| \geq \sigma n > \alpha n$  as required.

Now we prove that condition (b) of Lemma 7 holds. Consider a subproblem  $\mathcal{J}$  on  $\frac{1}{2}\alpha n \leq v \leq \alpha n$  variables and let H' be the hypergraph remaining after removing all cycle components from the underlying hypergraph of  $\mathcal{J}$ .

Case 1:  $|H'| \leq \frac{1}{4}\alpha n$ . Since  $|\mathcal{J}| \geq \frac{1}{2}\alpha n$ , the total size of the cycle components is at least  $\frac{1}{4}\alpha n$ . By Lemma 22 in Appendix B, a.s. at most  $\log n$  vertices lie on cycles of size at most 4. Every vertex on any other cycle component lies on a member of  $\mathcal{B}^2(\mathcal{J})$ , and each member of  $\mathcal{B}^2(\mathcal{J})$  contains fewer than 4k vertices. So  $|\mathcal{B}^2(\mathcal{J})| > \frac{1}{4k}(\frac{1}{4}\alpha n - \log n) > \frac{1}{20k}\alpha n > \zeta n$ .

Case 2:  $|H'| > \frac{1}{4}\alpha n$ . Since  $3\theta d(c/c^*)^{\theta} \leq \alpha/8$ , Lemma 16 yields that at least  $\alpha n/4 - \alpha n/8 > \alpha n/8$  of the variables in H' do not lie on any forcer paths of length at least  $\theta$ . Thus, any collection of forcer paths which covers all the variables of H' must contain at least  $\frac{\alpha}{8\theta}n$  paths. So we can apply Lemma 12 where  $\mathcal{I}$  is the CSP on H' formed by removing all cycle components from  $\mathcal{I}$  and where  $\frac{\alpha}{4\theta}n \leq q \leq |\mathcal{I}| \leq \alpha n$ . Since  $q \leq \alpha n < \sigma n$ , Lemma 15 implies that a.s. the entire random CSP has no  $Z_q$  configuration and so neither does H'. Since H' has no cycle components, Lemma 12 implies that  $\mathcal{B}(\mathcal{I}) \geq q/(72000k^3) \geq \zeta n$ . Every boundary element of  $\mathcal{I}$  is also a boundary element of  $\mathcal{I}$ , and so this establishes condition (b).

### 5 Proof of Theorem 3

We will show that a.s.  $\mathcal{F}_{n,M}^{d,k,t}$  contains a small unsatisfiable subproblem with a structure that is inspired by the snakes of [12]. Our proof is similar to the corresponding proof in [12].

A forbidding cycle is a  $x:\delta\to x':\delta'$  forcing path along with a  $x':\delta'\to x:\delta''$  forcer where  $\delta\neq\delta''$ . Thus, there is no satisfying assignment where  $x=\delta$ . We say that the cycle forbids  $x:\delta$ . Consecutive clauses in the cycle intersect in exactly one variable; such variables are called connecting variables.

An r-flower is the union of d forcing cycles  $C_1, \ldots, C_d$  such that: (i) each has exactly r forcers; (ii) each cycle contains a particular variable x; (iii) no other variable lies in more than one of

the cycles; (iv) each cycle  $C_i$  forbids x:i. We call x the center variable. Thus, any r-flower is unsatisfiable.

**Lemma 17** For any constants  $d, k \geq 2, d+k > 4$  and  $(d-1)d^{k-2} \leq t < d^{k-1}$ , and for every constant  $c > c^*(d, k, t)$ ,  $\mathcal{F}_{n, M=cn}^{d, k, t}$  a.s. contains all constraints of an r-flower, where  $r = \lambda \log n$ , for some sufficiently large constant  $\lambda$ .

Theorem 3 follows immediately since if  $\mathcal{F}_{n,M}^{d,k,t}$  contains an r-flower with  $r = \lambda \log n$ , then we can use exhaustive search to prove that the r-flower, and hence  $\mathcal{F}_{n,M}^{d,k,t}$ , is unsatisfiable in exp(r) = poly(n) steps. Such a proof can be simulated by a resolution proof with only a polynomial increase in length, so a.s.  $\mathbf{RES}(\mathcal{F}_{n,M}^{d,k,t}) = \text{poly}(n)$ . (The case k = d = 2, i.e. 2-SAT, is already known[12].)

**Proof:** For any potential flower A, we let  $X_A$  be the indicator variable for the event that the clauses of A all appear in the random CSP. With  $X = \sum_A X_A$ , it is enough to show that  $\mathbf{E}(X^2) \leq \mathbf{E}(X)^2(1+o(1))$ . Then the theorem follows easily from the Chebyshev inequality.

First, we compute  $\mathbf{E}(X)$ . We must choose s=dr-d+1 connecting variables (including the center variable) and the (k-2)dr other variables. There are s! ways to arrange the connecting variables and then  $[(k-2)dr]!/(k-2)!^{dr}$  ways to arrange the other variables into the flower. For each  $C_i$ , we need to choose some value other than i for the center variable and an arbitrary value for each of the other connecting variables. Then we multiply by the probability of all our forcers being present. We have  $c=(1+\epsilon)c^*$  for some  $\epsilon>0$ . For fixed variables x,y and values  $u,v\in D$ , the probability that there exists a  $(x:u)\longrightarrow (y:v)$  forcer is  $p=(1+\epsilon)\frac{(k-2)!}{dn^{k-1}}$ , by the same calculations as in Lemma 9.

$$\mathbf{E}(X) = \binom{n}{s} s! \binom{n-s}{(k-2)dr} \frac{[(k-2)dr]!}{(k-2)!^{dr}} \times (d-1)^d d^{(r-1)d} p^{dr}$$
$$= (1+o(1))(1+\epsilon)^{dr} \left(\frac{d-1}{d}\right)^r n^{1-d}$$
$$= \Omega(n^2)$$

for sufficiently large  $\lambda$ . Next we compute an upper bound for  $\mathbf{E}(X^2)$ . Consider a fixed r-flower A and its underlying hypergraph  $H_A$ . For each  $i, j \geq 1$  we will upper bound the number of potential r-flowers B that have exactly i constraints in common with A where these constraints form j connected components in  $H_A$ . (A very loose upper bound will suffice.) First, we consider choosing the j components. At most one component contains the center variable - for each  $C_i$ , such a component contains either all of  $C_i$ , none of  $C_i$ , or the portion of  $C_i$  between 2 variables. So there are at most  $(2+r^2)^d$  choices for such a component. Each of the other components is specified by 2 variables on the same cycle. Thus, the number of choices for the components of  $H_A$ is at most  $(2+r^2)^d(dr^2)^{j-1}$ . To obtain a very loose upper bound on the number of ways that these components can fit into B, we simply multiply by the number of ways to choose j components from B and then multiply by j! for the number of ways to pair them up with the components of A. Note that, since  $t < d^{k-1}$  and d, k are not both 2, no constraint can be a  $a: \delta \to b: \gamma$ forcer for more than one choice of  $a, b, \delta, \gamma$ . Therefore, if an edge lies in the underlying hypergraph of both A and B, then its forcer in B must be identical to its forcer in A. Therefore, to choose the rest of B, we choose the remaining at most s-i-j variables, a value for each of them (d values if one of them is the center variable) and then choose k-2 non-connecting variables for each of the remaining clauses. Thus, the total number of potential such r-flowers is at most  $((2+r^2)^d(dr^2)^{j-1})^2 j! n^{s-i-j} d^{s-i-j+d-1} \binom{n}{k-2}^{dr-i}$ . Therefore,  $\mathbf{E}(X^2)$  is:

$$\begin{split} & \sum_{A,B} \mathbf{Pr}(X_A = 1 \wedge X_B = 1) \\ & = \sum_{A} \sum_{B:A \cap B = \emptyset} \mathbf{Pr}(X_A = 1) \mathbf{Pr}(X_B = 1) + \sum_{A} \sum_{B:H_A \cap H_B \neq \emptyset} \mathbf{Pr}(X_A = 1 \wedge X_B = 1) \\ & < \mathbf{E}(X)^2 + \sum_{A} \mathbf{Pr}(X_A = 1) \times \sum_{i=1}^s \sum_{j=1}^i \left[ \left( (2 + r^2)^d (dr^2)^{j-1} \right)^2 \times j! n^{s-i-j} d^{s-i-j+d-1} \binom{n}{k-2}^{dr-i} p^{dr-i} \right] \\ & = \mathbf{E}(X)^2 + \mathbf{E}(X) \sum_{i=1}^s \left[ (2 + r^2)^{2d} r^{-4} n^{s-i} d^{s-i+d-3} \times \binom{n}{k-2}^{dr-i} p^{dr-i} \sum_{j=1}^i \left( \frac{dr^4 j}{n} \right)^j \right] \\ & \le \mathbf{E}(X)^2 + \mathbf{E}(X) \times \mathbf{E}(X) \times O(r^{4d-4}) \times \sum_{i=1}^s (1 + \epsilon)^{-i} O\left( \frac{r^4}{n} \right) \\ & \le \mathbf{E}(X)^2 \left( 1 + O\left( \frac{r^{4d}}{n} \right) \right), \end{split}$$

which completes the proof since  $r = O(\log n)$ .

### 6 Proofs of Theorems 4 and 5

We close the paper with the proofs of Theorems 4 and 5. Part (a) of Theorem 4 uses a very standard technique whereby we compute the expected number of satisfying assignments. There are other standard techniques around which will improve this theorem (eg. the techniques from [30]); we made no attempt to do so in part (a) as our only goal was to show that  $\mathcal{F}_{n,M=cn}^{d,k,t}$  is a.s. unsatisfiable for some  $c < c^*(d,k,t)$ . However, we need to work a bit harder for part (c) and so we use the simplest of the techniques from [30].

#### Proof of Theorem 4:

Part (a): Consider any instance  $\mathcal{I}$  chosen from  $\mathcal{F}_{n,M=cn}^{d,k,t}$ . There are  $d^n$  assignments to the variables of  $\mathcal{I}$ . For each such assignment, the probability that all constraints are satisfied is easily seen to be  $\left((d^k-t)/d^k\right)^{cn}$ . Therefore, the expected number of assignments that satisfy  $\mathcal{I}$  is

$$d^n \left(\frac{d^k - t}{d^k}\right)^{cn} = e^{(\ln d - c \ln(d^k/d^k - t))n},$$

which is o(1), if  $c > \frac{\ln d}{\ln(d^k/d^k-t)}$ . This implies that for such c, a.s.  $\mathcal{I}$  is unsatisfiable.

Part (b): It is straightforward to verify that the statement holds for d=2 and  $k \in \{4,5\}$ . Also,  $c^*(d,k,t) > \frac{1}{dk(k-1)} \times \left(\frac{d^k}{t}\right)^{(d-1)d^{k-2}} > \frac{d^{(d-1)d^{k-2}}}{dk(k-1)}$ . So it is enough to show that for  $d \ge 2$ ,  $k \ge 6$ , and for  $d \ge 3$ ,  $3 \le k \le 5$ :

$$\frac{\ln d}{\ln[d^k/(d^k-t)]} < \frac{d^{(d-1)d^{k-2}}}{dk(k-1)} \tag{1}$$

A simple inductive argument shows that for  $k \geq 6$ :  $k(k-1) < 2^k$  and  $k < 2^{k-2} - 4$ , which implies  $2k(k-1) < 2^{2^{k-2}-3} \leq d^{(d-1)d^{k-2}-3}$ . Similarly, for  $d \geq 3$  and  $k \geq 3$ :  $2k(k-1) < 3^k$  and  $k \leq 2 \times 3^{k-2} - 3$ , which implies  $2k(k-1) < 3^{2 \times 3^{k-2}-3} \leq d^{(d-1)d^{k-2}-3}$ . Therefore, in both cases:

$$\frac{2dk(k-1)\ln d}{d^{(d-1)d^{k-2}}} < \frac{d^{(d-1)d^{k-2}-2}\ln d}{d^2} < \frac{\ln d}{d^2} < \frac{(d-1)d^{k-2}}{d^k - (d-1)d^{k-2}} \tag{2}$$

Using (2) and the fact that for 0 < x < 1,  $e^x < 1 + 2x$ :

$$\exp\left(\frac{dk(k-1)\ln d}{d^{(d-1)d^{k-2}}}\right) \leq 1 + \frac{(d-1)d^{k-2}}{d^k - (d-1)d^{k-2}}$$

$$\leq \frac{d^k}{d^k - t}$$

$$\frac{dk(k-1)\ln d}{d^{(d-1)d^{k-2}}} < \ln\left(\frac{d^k}{d^k - t}\right)$$

$$\frac{\ln d}{\ln[d^k/(d^k - t)]} < \frac{d^{(d-1)d^{k-2}}}{dk(k-1)},$$

thus establishing (1) as required.

Part (c): For the case d=2, k=3, t=2, we strengthen the bound from part (a) by applying the so-called "1-flips" technique from [30]. We say that a satisfying assignment is 1-maximal if for every variable x with value 1, changing the value of x to 2 will result in a non-satisfying assignment. It is easy to see that if an instance  $\mathcal{I}$  chosen from  $\mathcal{F}_{n,M=cn}^{d,k,t}$  is satisfiable, then it has a 1-maximal satisfying assignment; indeed, consider an assignment which, amongst all satisfying assignments, has the greatest number of variables with value 2.

Consider any instance  $\mathcal{I}$  chosen from  $\mathcal{F}_{n,M=cn}^{d,k,t}$ . For each  $0 \leq a \leq n$ , consider one of the  $\binom{n}{a}$  value assignments  $\nu$  to the variables of  $\mathcal{I}$  in which exactly a variables have value 1. As described in the proof of part (a), the probability that this is a satisfying assignment is  $(6/8)^{cn}$ . We condition on the fact that it is a satisfying assignment; the effect of this conditioning is that the M constraints are chosen uniformly at random from amongst the  $\binom{7}{2}\binom{n}{3}$  constraints that are not violated by  $\nu$ .

For  $\nu$  to be 1-maximal, it must be the case that for each variable x with value  $\nu(x) = 1$ , there is a constraint on x and two other variables, say y, z that contains  $x = 2, y = \nu(y), z = \nu(z)$  as a restriction; we say that such a constraint blocks x. There are  $6\binom{n-1}{2}$  possible constraints of this form, 6 for each choice of y, z. Conditional on  $\nu$  being satisfying, the probability that no constraint blocks x is

$$\binom{\binom{7}{2}\binom{n}{3} - 6\binom{n-1}{2}}{M} / \binom{\binom{7}{2}\binom{n}{3}}{M} = e^{-6c/7} + o(1).$$

Given two variables,  $x_1, x_2$  with  $\nu(x_1) = \nu(x_2) = 1$ , the events that at least one constraint blocks  $x_1$  and at least one constraint blocks  $x_2$  are not independent - given that one blocks  $x_1$ , it is less likely that that constraint blocks  $x_2$  and so the probability that  $x_2$  is blocked is a bit smaller. However, it is easy to see that this dependence goes in the right direction for our purposes, and so the probability that all a variables that are assigned 1 by  $\nu$  are blocked is less than  $(1 - e^{-6c/7} + o(1))^a$ . Therefore, the expected number of 1-maximal satisfying assignments is

$$\sum_{a=0}^{n} \binom{n}{a} \left(\frac{6}{8}\right)^{cn} (1 - e^{-6c/7})^a = \left(\frac{6}{8}\right)^{cn} (2 - e^{-6c/7})^n = \left(\left(\frac{3}{4}\right)^c (2 - e^{-6c/7})\right)^n.$$

For c > 2.114 we have  $\left(\frac{3}{4}\right)^c (2 - e^{-6c/7}) < 1$  and so this expected number is o(1). Thus a.s.  $\mathcal{I}$  has no 1-maximal satisfying assignments and so  $\mathcal{I}$  is unsatisfiable.

(Clearly, this technique will easily give an improvement to the bound in part (a) for any d, k, t. It is not hard to see that further techniques from [30] will obtain even better bounds.)

Our final proof follows the, now rather standard, technique of using a differential equation analysis to show that a particular algorithm will a.s. find a satisfying assignment. See [1] for a good presentation of this method and survey of some of its most important applications.

**Proof of Theorem 5** Recall that our assumption is that  $c < c^*(2,3,3) = 7/9$  and that we wish to show that  $\mathcal{F}_{n,M=cn}^{2,3,3}$  is a.s. satisfiable.

We will make use of the fact that  $\mathcal{F}_{n,M=cn}^{2,3,3}$  has a sharp threshold in the sense of Friedgut's Theorem[22]. This fact was first proven independently by Creignou and Daudé[16] and by Istrate[27]. These papers each proved special cases of a more general conjecture from [15] which was proved in [17] where Creignou and Daudé classified which members of a large family of random CSP's with domain size two exhibit a sharp threshold.

More formally, this fact says that there is a function  $c^*(n)$  such that for every  $\epsilon > 0$ ,  $\mathcal{F}_{n,M=(c^*(n)-\epsilon)n}^{2,3,3}$  is a.s. satisfiable and  $\mathcal{F}_{n,M=(c^*+\epsilon)n}^{2,3,3}$  is a.s. unsatisfiable. This implies that if for some constant c,  $\mathcal{F}_{n,M=cn}^{2,3,3}$  is satisfiable with probability at least  $\gamma$  for some  $\gamma > 0$ , then  $\mathcal{F}_{n,M=c'n}^{2,3,3}$  is a.s. satisfiable for every constant c' < c: The only alternative is that  $\mathcal{F}_{n,M=cn}^{2,3,3}$  and  $\mathcal{F}_{n,M=c'n}^{2,3,3}$  are neither a.s satisfiable nor a.s. unsatisfiable. But this contradicts the existence of  $c^*(n)$  since for  $\epsilon = (c - c')/3$ , either  $c' < c^*(n) - \epsilon$  or  $c > c^*(n) + \epsilon$ .

Thus, it will suffice to prove that for every c < 7/9,  $\mathcal{F}_{n,M=cn}^{2,3,3}$  is satisfiable with probability at least  $\gamma$  for some  $\gamma = \gamma(c) > 0$ , which we do now.

We consider the following algorithm, which we denote Unit Constraint (UC).

The initial CSP is the input CSP, which is drawn from  $\mathcal{F}_{n,M=cn}^{2,3,3}$ . Repeatedly, we select a variable, x, and assign it a value i. We then modify each constraint C containing x as follows. If C contains any restrictions involving x:i then we form a new constraint C' on the variables of C other than x, by taking each restriction of C that contains x:i, removing x:i from that restriction, and placing the shortened restriction in C'. Thus, C' can be thought of as the constraint that is implied by C and setting x=i. Note that this might result in a constraint on exactly one variable in which each restriction simply dictates a value which that variable is not allowed to receive. We remove C and, if C contained any restrictions involving x:i, we replace it with C', unless:

- If C' is on two variables, say a, b, and if there is some value, j for which C' forbids both (a = 1, b = j) and (a = 2, b = j) then we simplify by replacing C with the constraint whose only variable is b and whose only restriction is (j), i.e. a 1-variable constraint that forbids b from taking the value j. If C' also forbids both (a = j', b = 1) and (a = j', b = 2) then we replace C by two 1-variable constraints which forbid a = j' and b = j, respectively. Note that this latter case occurs iff C' contains 3 restrictions (since C' cannot contain more than t = 3 restrictions).
- If C has exactly one variable, and if it forbids x = i then C' will contain no variables and so we remove C but do not add C', and we say that we formed a *null-constraint*. This indicates that our assignment violated one of the original constraints. However, we will continue to run the algorithm, as this will be convenient for its analysis.

Since no constraint has more than t = 3 restrictions, it is easy to see that we will never generate a 1-variable constraint with two restrictions.

We choose x and i as follows:

- If there are any clauses on 1 variable, choose one of them uniformly at random, and set it's variable so as to satisfy that clause. (As described above, no such clause will ever contain 2 restrictions and so there is always a unique such setting for that variable.)
- Otherwise, pick x uniformly at random from all unset variables and pick i uniformly at random from the domain  $\{1, 2\}$ .

After r variables have been set, we define the following:

- $C_3(r)$  the number of constraints on 3 variables
- $H_1(r)$  the number of constraints on 2 variables with 1 restriction
- $H_2(r)$  the number of constraints on 2 variables with 2 restrictions
- $C_1(r)$  the number of constraints on 1 variable

Note that the total number of remaining constraints is the sum of these values since every constraint formed has at least one restriction, and no constraint on 2 variables with 3 restrictions is ever added (since instead of adding such a constraint, we would add the equivalent pair of constraints each on 1 variable).

Claim 1: For each r, the CSP remaining after r steps of UC is uniformly random from amongst all CSP's with  $C_3(r)$  constraints with 3 variables and 3 restrictions,  $H_1(r)$  constraints with 2 variables and 1 restriction,  $H_2(r)$  constraints with 2 variables and 2 restrictions, and  $C_1(r)$  constraints with 1 variable and 1 restriction.

Proof: Consider two CSP's  $H_1$ ,  $H_2$  each with  $C_3(r)$  constraints with 3 variables and 3 restrictions,  $H_1(r)$  constraints with 2 variables and 1 restriction,  $H_2(r)$  constraints with 2 variables and 2 restrictions, and  $C_1(r)$  constraints with 1 variable and 1 restriction. Consider any input CSP  $F_1$  with n variables and with cn constraints each having 3 variables and 3 restrictions, and any sequence of random choices of variables, such that running UC with input  $F_1$  for r steps with that sequence of random choices will result in  $H_1$ . It is trivial to see how to modify  $F_1$  into  $F_2$ , also with n variables and with cn constraints each having 3 variables and 3 restrictions, such that running UC with input  $F_2$  for r steps with that same sequence of random choices will result in  $H_2$ . (Essentially, you simply replace all original constraints in  $F_1$  that became constraints in  $H_1$  with constraints that will instead become constraints of  $H_2$ .) This implies that the probability of ending up with  $H_1$  is the same as the probability of ending up with  $H_2$ . This, in turn, implies the claim.

Next, we consider the expected changes in the first 3 variables after step r+1. By examining all  $\binom{8}{3}$  possible constraints on 3 variables, it is straightforward (but tedious) to verify that we have the following, regardless of whether step r+1 sets the variable of a 1-clause or a uniformly random unset variable.

- $\mathbf{Exp}(C_3(r+1) C_3(r)) = -\frac{3C_3(r)}{n-r}$
- $\mathbf{Exp}(H_1(r+1) H_1(r)) = \frac{9}{7} \times \frac{C_3(r)}{n-r} \frac{2H_1(r)}{n-r}$
- $\mathbf{Exp}(H_2(r+1)) H_2(r)) = \frac{3}{7} \times \frac{C_3(r)}{n-r} \frac{2H_2(r)}{n-r}$

Furthermore, the expected number of new 1-variable constraints that are formed during step r+1 is:

$$F(r) = \frac{9}{7} \times \frac{C_3(r)}{n-r} + \frac{H_1(r)}{n-r} + \frac{2H_2(r)}{n-r}.$$

As is standard with this sort of analysis, our goal is to prove that a.s. F(r) is always less than  $1 - \zeta$  for some  $\zeta > 0$ , as this will imply that with sufficiently high probability, no null-constraints are formed; i.e., that a.s. our assignment does not violate any of the original constraints.

Consider the following functions:

- $c_3(x) = c(1-x)^3$ ,
- $h_1(x) = \frac{9c}{7}x(1-x)^2$
- $h_2(x) = \frac{3c}{7}x(1-x)^2$ .

Note that their derivatives satisfy

- $c_3'(x) = -\frac{3c_3}{1-x}$ ,
- $h'_1(x) = \frac{9}{7} \times \frac{c_3(x)}{1-x} \frac{2h_1(x)}{1-x}$
- $h_2'(x) = \frac{3}{7} \times \frac{c_3(x)}{1-x} \frac{2h_2(x)}{1-x}$ .

and that  $C_3(0) = c_1(0)n = n$ ,  $H_1(0) = h_1(0)n = 0$ ,  $H_2(0) = h_2(0)n = 0$ .

We also need to note that with very high probability, none of these parameters change much during any one iteration. In particular, it is straightforward to show that the probability of either  $C_3$ ,  $H_1$  or  $H_2$  changing by more than  $\log n$  during any one iteration is less than  $n^{-10}$ .

Noting the correspondence between these derivatives and the expected values computed above, Wormald's Theorem [37] implies that for any  $\alpha > 0$ , a.s. for every  $r \leq (1 - \alpha)n$  we have the following:

- $C_3(r) = c_3(r/n)n + o(n)$ ,
- $H_1(r) = h_1(r/n)n + o(n)$ ,
- $H_2(r) = h_2(r/n)n + o(n)$ .

(See [1] for a statement of Wormald's Theorem and a good discussion of how to apply it in settings like this one.)

Thus, a.s. for every  $r \leq (1-\alpha)n$ , setting x = r/n we have

$$F(r) = \frac{9}{7} \times \frac{c_3(x)}{1-x} + \frac{h_1(x)}{1-x} + \frac{2h_2(x)}{1-x} + o(1) = \frac{9c}{7}(1-x)(1+\frac{2}{3}x) + o(1) \le \frac{9c}{7} + o(1),$$

over the relevant range of  $0 \le x \le 1$ . Thus, a.s.  $F(r) < 1 - \zeta$  for some small constant  $\zeta > 0$ , so long as c < 7/9.

We will run UC until a point where at least  $\alpha n$  variables remain unset, for some particular small constant  $\alpha$  to be named later. First, we bound the probability of creating a null-constraint. The sequence  $C_1(r)$  follows the pattern of a random walk on the positive integers, with a barrier at 0, and with drift always bounded above by  $-\zeta$ . Standard arguments imply that a.s.  $\sum_{r=1}^{(1-\alpha)n} C_1(r) \leq Wn$  for some constant W. Note that, during step r, the probability that no null-constraint is formed is (if  $C_1(r) > 0$ ) equal to  $(1 - \frac{1}{2(n-r)})^{C_1(r)-1} > (1 - \frac{1}{2\alpha n})^{C_1(r)}$ . Therefore, the probability that no null-constraints are formed at all is at least  $(1 - \frac{1}{2\alpha n})^{Wn} = e^{-W/2\alpha} + o(1)$ .

Again, using the fact that  $C_1(r)$  has negative drift, it is straightforward to show that a.s. there is some  $(1-\alpha)n - \log^2 n < r \le (1-\alpha)n$  such that  $C_1(r) = 0$ . It will be convenient to halt UC there. With high probability, we are left with  $c_3(1-\alpha)n + o(n) < \alpha^3 n$  constraints on 3 variables

and  $(h_1(1-\alpha) + h_2(1-\alpha))n + o(n) < 2\alpha^2 n$  constraints on 2 variables. Let G be the hypergraph whose vertices are the unset variables, and where a group of vertices form a hyperedge iff they are the set of variables covered by a constraint.

It is straightforward to show that this random hypergraph G a.s. does not have a giant component. Perhaps the easiest way to see this is to add a third random vertex to each hyperedge of size 2, and call the resulting 3-uniform hypergraph G'. This leaves us with a random 3-uniform hypergraph on  $n' = \alpha n$  vertices and with fewer than  $3\alpha n' < n'/12$  hyperedges for  $\alpha < 1/36$ . It is well-known that the threshold for such a random hypergraph to have a giant component is when the number of hyperedges is n'/6, and in particular, that with probability at least  $\gamma'$  for some  $\gamma' > 0$ , every component of G' will be a tree. This would imply that every component of G is a tree. It is easy to see that if every component of G is a tree, then the formula is satisfiable.

Thus, the probability that the original formula is satisfiable is at least  $\gamma' \times e^{-W/2\alpha} + o(1) > 0$ , as required.

## Acknowledgments

We would like to thank Joe Culberson and Yong Gao for spotting an error in Theorem 4 of the conference version of this paper [35]. We also thank two anonymous referees for a very careful reading and many suggested improvements in the presentation. We are grateful for the support of NSERC, a Sloan Research Fellowship and a Premier's Research Excellence Award. Most of this research was completed while the first author was a visiting researcher at Microsoft Research and the second author was with Department of Computer Science, University of Toronto.

#### References

- [1] D. Achlioptas. Lower bounds for random 3-SAT via differential equations. Theoretical Computer Science, 265 (1-2), (2001), p.159-185.
- [2] D. Achlioptas, P. Beame and M. Molloy. A sharp threshold in proof complexity. J. Computer and System Sciences 68 (2), 238-268, (2004). Preliminary version in Proceedings of STOC 2001, 337 346.
- [3] D. Achlioptas, A. Chtcherba, G. Istrate and C. Moore. The phase transition in NAESAT and 1-in-k SAT. Proceedings of SODA 01, 721 722.
- [4] D. Achlioptas, L. Kirousis, E. Kranakis, and D. Krizanc Rigorous results for random (2 + p)-SAT. Proceedings of RALCOM '97 (1997), and Theoretical Computer Science, 265 (1-2), (2001), p.109-129.
- [5] D. Achlioptas, L. Kirousis, E. Kranakis, D. Krizanc, M. Molloy, and Y. Stamatiou. *Random constraint satisfaction: a more accurate picture*. Constraints **6**, 329 324 (2001). Preliminary version in Proceedings of CP 97, 107 120.
- [6] P. Beame, J. Culberson and D. Mitchell. The resolution complexity of random graph k-colorability. Discrete Applied Mathematics 153, 2005, 25 47.
- [7] P. Beame and T. Pitassi. Simplified and improved resolution lower bounds. Proceedings of FOCS 1996, 274 282.

- [8] P. Beame, R. Karp, T. Pitassi and M. Saks. *The efficiency of resolution and Davis-Putnam procedures*. SIAM Journal on Computing, **31**, 1048 1075 (2002). Preliminary version in Proceedings of STOC 1998.
- [9] E. Ben-Sasson and A. Wigderson. Short proofs are narrow resolution made simple. Journal of the ACM 48 (2001). Preliminary version in Proceedings of STOC 1999.
- [10] B. Bollobás. Random Graphs. Academic Press, London. (1985)
- [11] B. Bollobás, C. Borgs, J. T. Chayes, J. H. Kim, and D. B. Wilson. *The scaling window of the 2-SAT transition*. Random Structures and Algorithms **18** 201 256 (2001).
- [12] V. Chvátal and B. Reed. Mick gets some (the odds are on his side). In *Proceedings 33rd Annual Symposium on Foundations of Computer Science*, pages 620–627, Pittsburgh, PA, October 1992. IEEE.
- [13] V. Chvátal and E. Szemerédi. Many hard examples for resolution. Journal of the ACM 35 (1988) 759 768.
- [14] N. Creignou and H. Daude. Satisfiability threshold for random XOR-CNF formulas. Discrete Applied Mathematics, **96-97**, 4153 (1999)
- [15] N. Creignou and H. Daudé. Generalized satisfiability problems: minimal elements and phase transitions. Theor. Comput. Sci. 1-3(302), 417-430 (2003). Preliminary version: Random generalized satisfiability problems. Proceedings of SAT 2002.
- [16] N. Creignou and H. Daudé. Combinatorial sharpness criterion and phase transition classification for random CSPs. Inf. Comput. 190(2), 220-238 (2004). Preliminary version: Crossed classification, complexity versus phase transition, for Boolean CSPs. Proceedings of SAT 2003.
- [17] N. Creignou and H. Daudé. Coarse and sharp transitions for random generalized satisfiability problems. In Mathematics and Computer Science III: Algorithms, Trees, Combinatorics and Probabilities M. Drmota, P. Flajolet, D. Gardy and B. Gittenberger (Eds). Birkhauser-Springer (2004) 507 - 516. Proceedings of the Third Colloquium on Mathematics and Computer Science, Vienna, 2004.
- [18] N. Creignou, H. Daude and O. Dubois. Approximating the satisfiability threshold of random k-XOR-formulas. Combinatorics, Probability, and Computing 12(2), 2003, 113 126.
- [19] O. Dubois and J. Mandler. The 3-XORSAT threshold (strong evidence in favour of the replica method of statistical physics). Proceedings of FOCS 2002.
- [20] M. Dyer, A. Frieze and M. Molloy. A probabilistic analysis of randomly generated binary constraint satisfaction problems. Theoretical Computer Science 290(3), 2003, 1815 1828.
- [21] W. Fernandex de la Vega. On random 2-SAT. manuscript (1992).
- [22] E. Friedgut and an appendix by J. Bourgain. Sharp thresholds of graph properties and the k-SAT problem. J. American Math. Soc. 12 (1999), 1017 1054.
- [23] Y. Gao and J. Culberson. Resolution complexity of random constraint satisfaction problems: another half of the story. Discrete Applied Mathematics 153, 2005, 124 - 140. Preliminary version in Proceedings of LICS'03, Workshop on Typical Case Complexity and Phase Transitions.

- [24] I. Gent, E. MacIntyre, P. Prosser, B. Smith and T. Walsh. Random constraint satisfaction: flaws and structure. Constraints 6, 345 372 (2001).
- [25] A. Goerdt. A threshold for unsatisfiability. Journal of Computer and System Sciences, 53:469–486, 1996.
- [26] A. Haken. The intractability of resolution. Theoretical Computer Science 39, 297 305 (1985).
- [27] G. Istrate. Threshold properties of random boolean constraint satisfaction problems. Disc. Appl. Math. 153, 141 - 152 (2005). Preliminary version in Proceedings of LICS'03, Workshop on Typical Case Complexity and Phase Transitions.
- [28] R. Karp. The transitive closure of a random digraph. Random Structures and Algorithms 1 73 93 (1990).
- [29] S. Janson, T. Łuczak and A. Ruciński. Random Graphs. Wiley, New York (2000).
- [30] L. Kirousis, E. Kranakis, D. Krizanc and Y. Stamatiou. Approximating the unsatisfiability threshold of random formulas. Random Structures and Algorithms 12 (1998), 253 269.
- [31] T. Łuczak. Size and connectivity of the k-core of a random graph. Discrete Math. 91, 61 68 (1991).
- [32] D. Mitchell. Resolution complexity of random constraints. Proceedings of Principles and Practices of Constraint Programming CP 2002.
- [33] D. Mitchell *The Resolution Complexity of Constraint Satisfaction*. Ph.D. Thesis, University of Toronto, 2002.
- [34] M. Molloy, Models for Random Constraint Satisfaction Problems. SIAM Journal on Computing 32(4), 2003, 935 949. Preliminary version in Proceedings of STOC 2002, 209 217.
- [35] M. Molloy and M. Salavatipour. The resolution complexity of random constraint satisfaction problems. Proceedings of FOCS 2003.
- [36] R. Monasson and R. Zecchina. Tricritical point in the random 2+p-SAT problem. J. Phys. A 31, 9209 (1998).
- [37] N. Wormald. Differential equations for random processes and random graphs. Ann. Appl. Prob. 5:4 1217 1235 (1995).
- [38] R. Monasson, R. Zecchina, S. Kirkpatrick, B. Selman and L. Troyansky. 2+p-SAT: Relation of typical-case complexity to the nature of the phase transition. Random Structure and Algorithms 15, 414 (1999).
- [39] R. Monasson, R. Zecchina, S. Kirkpatrick, B. Selman and L. Troyansky. *Phase transition and search cost in the 2+p-sat problem*. Proceedings of PhysComp 96, T. Toffoli, M. Biafore, J. Leao eds., Boston (1996).
- [40] A. Urquhart. Hard examples for resolution. Journal of the ACM, 34 209 219 (1987).
- [41] K. Xu and W. Li. Exact Phase transitions in random constraint satisfaction problems. J. Artificial Intelligence Research, 12(2000):93-103.

## Appendix A: a short proof of Theorem 6

Each variable v has two literals: v and  $\overline{v}$ , and we say that they are *complements* of each other. A literal of a variable in a CNF-formula F, is *pure* if does not appear in any clause of F. As explained in [2], it suffices to prove:

**Lemma 18** For any  $\Delta$ ,  $\epsilon > 0$ , consider a random CNF-formula  $\mathcal{F}$  on n variables with  $\Delta n$  3-clauses and  $(1 - \epsilon)n$  2-clauses where every such formula is equally likely. A.s.:

- (a) every subformula on at most αn variables is satisfiable, and
- (b) every subformula on v variables where  $\frac{1}{2}\alpha n \leq v \leq \alpha n$  has at least  $\zeta n$  pure literals.

Consider any CNF-formula F containing only 2-clauses and 3-clauses. Let  $H_2$  be the graph defined as follows: the vertices of  $H_2$  are the variables of F, and two vertices are joined iff their variables form at least one 2-clause.

An isolated cycle of F is a 2-SAT subformula F' such that (i) F' forms a component of  $H_2$  that is a cycle, (ii) no variable of F' lies in a 3-clause and (iii) no variable of F' has a pure literal.

A pendant path of F is a path of  $H_2$  whose internal vertices each have degree 2 in  $H_2$  and do not lie in any 3-clauses. Consider such a path whose variables are, in order:  $x_0, x_1, ..., x_l$ , and suppose that  $C_1, ..., C_l$  is the corresponding set of 2-clauses in F. If, for each  $x_i$ , the literal of  $x_i$  appearing in  $C_i$  is the complement of the literal appearing in  $C_{i+1}$ , then we say that it is a forcing path. Note that such a path is bidirectional in the sense that there are two literals a, b such that in any satisfying assignment, if  $x_0 = a$  then  $x_l = b$ , and if  $x_l = \overline{b}$  then  $x_0 = \overline{a}$ . Trivially, a single vertex is a forcing path of length 0.

For any  $r \geq 1$ , a  $Y_r$ -configuration consists of:

- r forcing paths; and
- a collection of  $t_2$  additional 2-clauses and  $t_3$  3-clauses whose variables are all endpoints of the r forcing paths, for some  $t_2, t_3$  with  $\frac{3}{2}t_2 + 3t_3 \ge \frac{5}{3}r$ .

Consider a collection of forcing paths  $\mathcal{P} = P_1, ..., P_r$  of F such that (i) every variable of F appears on exactly one path and (ii)  $\mathcal{P}$  is minimal in the sense that it is impossible to form a collection  $P_1, ..., P_{r-1}$  satisfying (i) by adding a 2-clause from F to  $\mathcal{P}$ . Obviously a collection of paths of length 0, one for each variable of F satisfies (i), and so some collection exists which satisfies (i) and (ii).

**Lemma 19** If F has at most r/3 pure literals and no isolated cycles, then F has a  $Y_r$ -configuration.

**Proof:** We call the clauses of  $\mathcal{P}$  path clauses and the other clauses in F non-path clauses. Note that every non-path clause only contains variables that are endpoints of the paths in  $\mathcal{P}$ . We define a set X of literals as follows: For each  $P_i$  of length 0, we place both literals of the variable of  $P_i$  into X. For each  $P_i$  of length at least one, and for each endpoint v of  $P_i$ , we place the literal of v that does not appear on a clause of  $P_i$  into X. Thus, |X| = 2r, and any literal of X that does not appear in a non-path clause is pure.

We form a graph G with vertex set X as follows: Every non-path 2-clause  $(a \lor b)$  forms an edge of G: x(a) is defined to be a if  $a \in X$  and the complement of a otherwise; x(b) is defined in the same way; the edge of G is between x(a) and x(b).

Let s be the number of pure literals in F,  $t_2 = |E(G)|$  be the number of 2-clauses and  $t_3$  be the number of 3-clauses.

Claim: Every component of G with either 1 or 2 literals contains a literal which is either pure or in a 3-clause.

*Proof:* If the component has 1 literal, then it is either pure or in a 3-clause. So suppose the component has 2 literals and the edge joining them corresponds to the clause  $(a \lor b)$ . If at least one of a, b is not in X, then the complement of that literal is either pure or in a 3-clause. If a and b are both in X and if their variables are the endpoints of the same path in  $\mathcal{P}$ , then one of them must be in a 3-clause, or else that path plus  $(a \lor b)$  would form an isolated cycle. If a and b are both in X and their variables are the endpoints of different paths of  $\mathcal{P}$ , then it is easy to see that those two paths plus  $(a \lor b)$  form a forcing path. This contradicts the minimality of  $\mathcal{P}$ , i.e. condition (ii) in the definition of  $\mathcal{P}$ .

Let  $\ell_1$  be the number of components with exactly one literal, and  $\ell_2$  be the number with exactly 2 literals. The remaining components have  $2r - \ell_1 - 2\ell_2$  literals and at least  $\frac{2}{3}(2r - \ell_1 - 2\ell_2)$  edges (since the components on those literals each have size at least 3). Therefore,  $t_2 \geq \ell_2 + \frac{2}{3}(2r - \ell_1 - 2\ell_2)$ ; i.e.  $\frac{3}{2}t_2 \geq 2r - \ell_1 - \frac{1}{2}\ell_2$ . By Claim 1,  $3t_3 + s \geq \ell_1 + \ell_2$ . These combine to yield:

$$\frac{3}{2}t_2 + 3t_3 + s \ge 2r.$$

Since  $s \leq r/3$ , F has a  $Y_r$  configuration.

**Lemma 20** For any  $\Delta$ ,  $\epsilon > 0$ , consider a random CNF-formula  $\mathcal{F}$  on n variables with  $\Delta n$  3-clauses and  $(1 - \epsilon)n$  2-clauses where every such formula is equally likely. There is some constant  $\alpha > 0$  such that a.s. F has no  $Y_r$  configuration for any  $r \leq \alpha n$ .

**Proof:** Given  $r, t_2$  we specify  $t_3 = \lceil \frac{5}{9}r - \frac{1}{2}t_2 \rceil$  to be the smallest integer  $t_3$  such that  $\frac{3}{2}t_2 + 3t_3 \ge \frac{5}{3}r$ . Clearly, it suffices to show that a.s. there are no  $Y_r$  configurations with such a pair  $t_2, t_3$ . First, we compute the expected number of  $Y_r$  configurations for any choice of  $t_2, t_3$  that are both at least r/100.

Consider any list of 2-clauses  $C_1, ..., C_s$ . The probability that they all appear in  $\mathcal{F}$  is

$$\frac{\binom{4\binom{n}{2}-s}{(1-\epsilon)n-s}}{\binom{4\binom{n}{2}}{(1-\epsilon)n}} < \left(\frac{1-\epsilon}{2(n-1)}\right)^s < \left(\frac{1-\epsilon'}{2n}\right)^s$$

for some  $0 < \epsilon' < \epsilon$ .

We have at most  $\binom{n}{r}n^r$  choices for the r pairs of endpoints. Suppose that the numbers of 2-clauses in the paths are  $l_1, ..., l_r$ , and set  $L = l_1 + ... + l_r$ . Then there are  $n^{L-r}$  choices for the interior variables on the paths, and  $2^{L+r}$  choices for the literals. We multiply by the probability that all L of these clauses appear and that there are  $t_2$  other 2-clauses and  $t_3$  3-clauses on the endpoints. This gives us an upper bound of

$$\sum_{l_1,\dots,l_r\geq 0} \binom{n}{r} n^L 2^{L+r} \left(\frac{1-\epsilon'}{2n}\right)^L \binom{(1-\epsilon)n}{t_2} \binom{\Delta n}{t_3} \left(\frac{2r}{n}\right)^{2t_2+3t_3}$$

$$\leq \left(\frac{2en}{r}\right)^r \left(\frac{en}{t_2}\right)^{t_2} \left(\frac{e\Delta n}{t_3}\right)^{t_3} \left(\frac{2r}{n}\right)^{2t_2+3t_3} \sum_{l_1,\dots,l_r>0} (1-\epsilon')^L$$

$$\leq \left(\frac{\gamma r}{n}\right)^{t_2 + 2t_3 - r} \left(\sum_{l \geq 0} (1 - \epsilon')^l\right)^r \quad \text{for some } \gamma > 0, \text{ since } t_2, t_3 \geq r/100$$

$$\leq \left(\frac{\gamma' r}{n}\right)^{r/9}$$

For some constant  $\gamma' > 0$  since  $t_2 + 2t_3 - r = \frac{2}{3}(\frac{3}{2}t_2 + 3t_3) - r \ge \frac{2}{3} \times \frac{5}{3}r - r = \frac{r}{9}$ . If  $t_2 \le r/100$  then  $t_3 \ge (\frac{5}{9} - \frac{1}{200})r$ . For such  $t_2$ , we compute the expected number of collections of r vertex disjoint forcing paths along with  $t_3$  3-clauses on their endpoints. Clearly, if there are no such collections then there is no  $Y_r$  configuration with those values of  $t_2, t_3$ . As above, we upper bound this expected number with:

$$\left(\frac{2en}{r}\right)^r \left(\frac{e\Delta n}{t_3}\right)^{t_3} \left(\frac{2r}{n}\right)^{3t_3} \left(\sum_{l\geq 0} (1-\epsilon')^l\right)^r < \left(\frac{\gamma'r}{n}\right)^{r/10},$$

with possibly an increase in  $\gamma'$ . If  $t_3 \leq r/100$  then we compute the expected number of collections of r vertex disjoint forcing paths along with t<sub>2</sub> additional 2-clauses on their endpoints. Clearly, if there are no such collections then there is no  $Y_r$  configuration with those values of  $t_2, t_3$ . Again, this expected number is at most  $\left(\frac{\gamma' r}{n}\right)^{r/10}$ . Thus, considering that there are O(r) choices for  $t_2, t_3, t_4$ it suffices to show that:

$$\sum_{r=1}^{\alpha n} r \left(\frac{\gamma' r}{n}\right)^{r/10} = o(1).$$

The first  $\log n$  terms of this sum add up to at most  $O(\log n/n^{1/10})$  and if  $\alpha < \frac{1}{2\gamma'}$  then the rest add up to at most  $\sum_{i > \log n} (1/2)^i = o(1)$ .

**Lemma 21** A.s. our random  $\mathcal{F}$  has at most log n variables lying on isolated cycles.

**Proof:** Consider the subformula  $F_2$  formed by the 2-clauses of  $\mathcal{F}$ . For any variable v lying in an isolated cycle of  $\mathcal{F}$ , there must be a sequence of 2-clauses in  $F_2$  of the form:  $(v \vee x_1), (\overline{x_1} \vee x_2)$  $(x_2), \dots, (\overline{x_i} \vee \overline{v})$ . A well-known property of random 2-SAT (see eg [11]) says that a.s. there are at most  $\log n$  such variables.

**Proof of Lemma 18:** If F' is a minimally unsatisfiable subformula of  $\mathcal{F}$ , then F' must be connected and F' cannot be an isolated cycle. Therefore F' can have no isolated cycles. Furthermore F' can have no pure literals. Therefore, by Lemma 19, F' must have a  $Y_r$  configuration for some  $r \geq 1$ . Therefore, by Lemma 20,  $\mathcal{F}$  a.s. has no minimally unsatisfiable subformula on at most  $\alpha n$ variables and hence a.s. has no unsatisfiable subformula on at most  $\alpha n$  variables. This establishes part (a).

Consider any subformula on v variables where  $\frac{1}{2}\alpha n \leq v \leq \alpha n$ . By Lemma 21,  $\mathcal{F}$  is a.s. such that after removing all isolated cycles from such a subformula, we are left with a subformula  $F^\prime$  on v' variables where  $\frac{1}{2}\alpha n - \log n \le v' \le \alpha n$ . It will suffice to show that a.s. every such F' has at least (n) pure literals. By Lemma 20,  $\mathcal{F}$  is a.s. such that every subformula on at most  $\alpha n$  variables does not have a  $Y_r$  configuration for any  $r \geq 1$ . Therefore, by Lemma 19, there is some  $r \geq 1$  such that F' has at least r/3 pure literals and F' has a collection  $\mathcal{P}$  of r forcer paths which contain all of its variables.

Well-known properties of random 2-SAT (see eg [11]) imply that there is some  $\pi > 0$  such that for every  $\theta > 0$ , a.s. F has fewer than  $e^{-\pi\theta}$  variables that lie on forcing paths of length at least  $\theta$ . In fact, the same would be true if we removed from the definition of forcing path the stipulation that no internal variables lie in any 3-clauses, thus reducing forcing paths to be defined only in terms of the  $(1 - \epsilon)n$  random 2-clauses of F. Pick  $\theta$  so that  $e^{-\pi\theta} < \alpha/4$ . Thus, at least  $(\alpha/4)n - \log n$  variables of F' lie on paths in  $\mathcal{P}$  of length less than  $\theta$ . Therefore,  $r > \alpha n/(5\theta)$  and so F' has at least  $\zeta n$  pure literals for  $\zeta = \alpha/(15\theta)$ . This establishes part (b).

## Appendix B: two simple lemmas

Here, we translate two standard facts from random graph theory into the setting of this paper.

**Lemma 22** A.s. the underlying random hypergraph of  $\mathcal{F}_{n,M}^{d,k,t}$  has fewer than  $\log n$  cycles of length at most 4.

**Remark:** This statement remains true when "4" is replaced by any constant. It is well known for random graphs, but we can't find the statement for random hypergraphs recorded in the literature. So we include the simple proof here.

**Proof:** For each constant integer  $r \geq 2$ , we compute the expected number of cycles of length r, by pretending that the vertices of degree 2 in the cycle are labelled  $v_1..., v_r$ , in order around the cycle; this creates an overcount which we correct by dividing by 2r, the number of ways to label each cycle. (For k = 2, we take  $r \geq 3$  since a simple graph contains no 2-cycles; for r = 2 we must divide by r instead of 2r since there are only 2 ways to label each cycle.)

There are at most  $n^r$  choices for  $v_1, ..., v_r$ , and at most  $\binom{n-r}{k-2}^r$  choices for the other vertices. This specifies the r edges. There are  $M^r = (cn)^r$  choices for which random edges correspond to the edges of the cycle. The probability that each such random edge is the desired one is  $\binom{n}{k}^{-r}$ . So the expected number of cycles of length r is:

$$\frac{1}{2r}n^r \binom{n-r}{k-2}^r (cn)^r \binom{n}{k}^{-r} = O(1).$$

(For r=2 we replace  $\frac{1}{2r}$  by  $\frac{1}{r}$ .) Thus, the expected number of cycles of length  $2 \le r \le 4$  is O(1). So by Markov's Inequality, the probability that this number is at least  $\log n$  is  $O(1/\log n) = o(1)$  as required.

**Remark:** With more work, one can show that the probability is much lower than  $O(1/\log n)$ ; but that is not needed here.

Our second lemma shows how the  $\mathcal{F}_{n,M}^{d,k,t}$  and  $\mathcal{F}_{n,p}^{d,k,t}$  models are, in many senses, equivalent and in particular, allows us to use the  $\mathcal{F}_{n,p}^{d,k,t}$  model in the proofs of Lemmas 15 and 16.

We say that a property A of CSP's in  $\Omega^{d,k,t}$  is monotone increasing if for every  $F_1, F_2 \in \Omega^{d,k,t}$  with every constraint of  $F_1$  also in  $F_2$ , if  $F_1$  has A then so does  $F_2$ . A is monotone decreasing if the same holds whenever every constraint of  $F_2$  is also in  $F_1$ . A is monotone if it is either montone increasing or monotone decreasing.

For example, it is easy to see that the properties considered in the statements of Lemmas 15 and 16 are both monotone decreasing.

**Lemma 23** Let A be any monotone property of CSP's in  $\Omega^{d,k,t}$ , and let c > 0 be any positive constant. A holds a.s. for  $\mathcal{F}_{n,p}^{d,k,t}$  with  $p = c \times k!/n^{k-1}$  iff for every real constant x, A holds a.s. for  $\mathcal{F}_{n,M}^{d,k,t}$  with  $M = \lceil cn + x\sqrt{n} \rceil$ .

In particular, taking x = 0 allows us to show that A holds a.s. in  $\mathcal{F}_{n,M}^{d,k,t}$  by proving that it holds a.s. in  $\mathcal{F}_{n,p}^{d,k,t}$ , as we do in the proofs of Lemmas 15 and 16. The proof is very straightforward, and follows similar proofs in, eg. [10, 29].

**Proof:** We assume that A is monotone increasing (the monotone decreasing case is nearly identical).

Suppose that  $\mathcal{F}_{n,p}^{d,k,t}$  a.s. has A. Fix any real x and set  $M = \lceil cn + x\sqrt{n} \rceil$ . Let  $\gamma(n) = \mathbf{Pr}(\mathcal{F}_{n,M}^{d,k,t}$  does not have A). The probability that the number of edges in  $\mathcal{F}_{n,p}^{d,k,t}$  is at most M is well known to be at least a positive constant g(x), since this number is a binomial variable with mean cn. By the monotonicity of A, the probability that  $\mathcal{F}_{n,p}^{d,k,t}$  does not have A is at least  $g(x) \times \gamma(n)$ . Therefore,  $\lim_{n \to \infty} \gamma(n) = 0$ , as required.

For the other direction, suppose that for every real constant x, A holds a.s. for  $\mathcal{F}_{n,M}^{d,k,t}$  with  $M = \lceil cn + x\sqrt{n} \rceil$ . For any  $\epsilon > 0$ , there exists  $x_1 < 0 < x_2$  such that the probability that the number of edges in  $\mathcal{F}_{n,p}^{d,k,t}$  is in  $(\lceil cn + x_1\sqrt{n} \rceil, ..., \lceil cn + x_2\sqrt{n} \rceil)$  is at least  $1 - \epsilon$ . Therefore, the probability that  $\mathcal{F}_{n,p}^{d,k,t}$  does not have A is at most  $\epsilon + o(1)$ . Since this is true for every  $\epsilon > 0$ ,  $\mathcal{F}_{n,p}^{d,k,t}$  a.s. has A. (Note that this part did not require A to be monotonic.)