

The Satisfiability Threshold for Non-Uniform Random 2-SAT*

Tobias Friedrich¹ and Ralf Rothenberger¹

¹Hasso Plattner Institute, University of Potsdam, Potsdam, Germany, `firstname.lastname@hpi.de`

Abstract

Propositional satisfiability (SAT) is one of the most fundamental problems in computer science. Its worst-case hardness lies at the core of computational complexity theory, for example in the form of NP-hardness and the (Strong) Exponential Time Hypothesis. In practice however, SAT instances can often be solved efficiently. This contradicting behavior has spawned interest in the average-case analysis of SAT and has triggered the development of sophisticated rigorous and non-rigorous techniques for analyzing random structures.

Despite a long line of research and substantial progress, most theoretical work on random SAT assumes a *uniform* distribution on the variables. In contrast, real-world instances often exhibit large fluctuations in variable occurrence. This can be modeled by a *non-uniform* distribution of the variables, which can result in distributions closer to industrial SAT instances.

We study satisfiability thresholds of non-uniform random 2-SAT with n variables and m clauses and with an arbitrary probability distribution $(p_i)_{i \in [n]}$ with $p_1 \geq p_2 \geq \dots \geq p_n > 0$ over the n variables. We show for $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ that the asymptotic satisfiability threshold is at $m = \Theta\left(\frac{1 - \sum_{i=1}^n p_i^2}{p_1 \cdot (\sum_{i=2}^n p_i^2)^{1/2}}\right)$ and that it is coarse. For $p_1^2 = o(\sum_{i=1}^n p_i^2)$ we show that there is a sharp satisfiability threshold at $m = (\sum_{i=1}^n p_i^2)^{-1}$. This result generalizes the seminal works by Chvatal and Reed [FOCS 1992] and by Goerdts [JCSS 1996].

*This paper is partially funded by the project *Skalenfreie Erfüllbarkeit* (project no. 416061626) of the German Research Foundation (DFG).

1 Introduction

Satisfiability of Propositional Formulas (SAT) is one of the most thoroughly researched topics in theoretical computer science. It was one of the first problems shown to be NP-complete by Cook [15] and, independently, by Levin [30]. Today SAT stands at the core of many techniques in modern complexity theory, for example NP-completeness proofs [29] or running time lower bounds assuming the (Strong) Exponential Time Hypothesis [10, 17, 26, 27].

In addition to its importance for theoretical research, Propositional Satisfiability is also famously applied in practice. Despite the theoretical hardness of SAT, many problems arising in practice can be transformed to SAT instances and then solved efficiently with state-of-the-art solvers. Problems like hard- and software verification, automated planning, and circuit design are often transformed into SAT instances. Such formulas arising from practical and industrial problems are therefore referred to as *industrial SAT instances*. The efficiency of SAT solvers on these instances suggests that they have a structure that makes them easier to solve than the theoretical worst-case.

1.1 Uniform Random k -SAT and the satisfiability threshold conjecture:

Random k -SAT is used to study the average-case complexity of Boolean Satisfiability. In the model, a random formula Φ with n variables, m clauses, and k literals per clause is generated in conjunctive normal form. Each of these formulas has the same uniform probability to be generated. Therefore, we also refer to this model as *uniform random k -SAT*.

One of the most prominent questions related to studying uniform random k -SAT is trying to prove the *satisfiability threshold conjecture*. The conjecture states that for a uniform random k -SAT formula Φ with n variables and m clauses there is a real number r_k such that

$$\lim_{n \rightarrow \infty} \Pr(\Phi \text{ is satisfiable}) = \begin{cases} 1 & m/n < r_k; \\ 0 & m/n > r_k. \end{cases}$$

Chvatal and Reed [11] and, independently, Goerdt [24] proved the conjecture for $k = 2$ and showed that $r_2 = 1$. For larger values of k upper and lower bounds have been established, e. g., $3.52 \leq r_3 \leq 4.4898$ [18, 25, 28]. Methods from statistical mechanics [32] were used to derive a numerical estimate of $r_3 \approx 4.26$. Coja-Oghlan and Panagiotou [12, 13] showed a bound (up to lower order terms) for $k \geq 3$ with $r_k = 2^k \log 2 - \frac{1}{2}(1 + \log 2) \pm o_k(1)$. Finally, Ding, Sly, and Sun [19] proved the exact position of the threshold for sufficiently large values of k . Still, for k between 3 and the values determined by Ding, Sly, and Sun the conjecture remains open.

The satisfiability threshold is also connected to the average hardness of solving instances. For uniform random k -SAT for example, the on average hardest instances are concentrated around the threshold [33].

1.2 Non-Uniform Random SAT

There is a large body of work which considers other random SAT models, e. g. regular random k -SAT [7, 8, 14, 39], random geometric k -SAT [9] and $2 + p$ -SAT [1, 34–36]. However, most of these are not motivated by modeling the properties of industrial instances. One such property is community structure [4], i. e. some variables have a bias towards appearing together in clauses. It is clear by definition that such a bias does not exist in uniform random k -SAT. Therefore, Giráldez-Cru and Levy [23] proposed the Community Attachment Model, which creates random formulas with clear community structure. However, the work of Mull et al. [38] shows that instances generated by this model have exponentially long resolution proofs with high probability, making them hard on average for solvers based on conflict-driven clause learning.

Another important property of industrial instances is their degree distribution. The degree distribution of a formula Φ is a function $f: \mathbb{N} \rightarrow \mathbb{N}$, where $f(x)$ denotes the fraction of different Boolean variables that appear x times in Φ (negated or unnegated). Instances created with the uniform random k -SAT model have a binomial distribution, while some families of industrial instances appear to follow a power-law distribution [2], i. e. $f(x) \sim x^{-\beta}$, where β is a constant intrinsic to the instance. Therefore, Ansótegui et al. [3] proposed a random k -SAT model with a power-law degree distribution. Empirical studies by the same authors [2, 3, 5, 6] found that

this distribution is beneficial for the runtime of SAT solvers specialized in industrial instances. However, it looks like instances generated with their model can be solved faster than uniform instances, but not as fast as industrial ones: median runtimes around the threshold still seem to scale exponentially for several state-of-the-art solvers [22].

Therefore, we want to consider a generalization of the model by Ansótegui et al. [2]. Our model allows instances with *any* given ensemble of variable distributions instead of only power laws: We draw m clauses of length k at random. For each clause the k variables are drawn with a probability proportional to the n -th distribution in the ensemble, then they are negated independently with a probability of $1/2$ each. This means, the probability ensemble is part of the model, but the number of variables n determines which distribution from the ensemble we actually use. We call this model *non-uniform random k -SAT* and denote it by $\mathcal{D}(n, k, (\vec{p}_x)_{x \in \mathbb{N}}, m)$. Although $\mathcal{D}(n, k, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ cannot capture all properties of industrial instances, e.g. community structure, it can help us to investigate the influence of the degree distribution on the structure and on the computational complexity of such instances in an average-case scenario.

As one of the steps in analyzing this connection, we would like to find out for which ensembles of variable probability distributions an equivalent of the satisfiability threshold conjecture holds in non-uniform random k -SAT. In previous works we already proved upper and lower bounds on the threshold position [21] and showed sufficient conditions on sharpness [20]. In this work we are interested in actually determining the satisfiability threshold for $k = 2$. It has to be noted that Cooper et al. [16] and Levy [31] already studied thresholds in a similar random 2-SAT model. The difference is that in their models the degrees are fixed and the random instances determined in a configuration-model-like fashion, while in our model we only have a sequence of expected degrees from which the actual degrees might deviate. Another difference is that we do a complete analysis of the model we consider, while they have additional constraints on their degree sequences. However, if we assume the expected degrees that our model implies to be the actual degrees, the thresholds determined by Cooper et al. and by Levy coincide with the ones we derive for our model.

1.3 Our Results

We investigate the position and behavior of the satisfiability threshold for non-uniform random 2-SAT. That is, we fix the number of variables n and the variable distribution \vec{p}_n from the ensemble and vary the number of clauses $m(n)$. To this end, we use the following definition and say that a function $m^*(n)$ is an *asymptotic satisfiability threshold* if

$$\Pr_{\Phi \sim \mathcal{D}(n, k, (\vec{p}_x)_{x \in \mathbb{N}}, m)}(\Phi \text{ satisfiable}) = \begin{cases} 1 - o(1) & \text{if } m(n) = o(m^*(n)) \\ o(1) & \text{if } m(n) = \omega(m^*(n)). \end{cases}$$

We also say that an asymptotic satisfiability threshold $m^*(n)$ is *sharp* if for all $\varepsilon > 0$

$$\Pr_{\Phi \sim \mathcal{D}(n, k, (\vec{p}_x)_{x \in \mathbb{N}}, m)}(\Phi \text{ satisfiable}) = \begin{cases} 1 - o(1) & \text{if } m(n) \leq (1 - \varepsilon) \cdot m^*(n) \\ o(1) & \text{if } m(n) \geq (1 + \varepsilon) \cdot m^*(n). \end{cases}$$

If an asymptotic threshold is not sharp, we call it *coarse*.

Let $\vec{p}_n = (p_1, p_2, \dots, p_n)$ be the variable probability distribution we use. W.l.o.g. we assume $p_1 \geq p_2 \geq \dots \geq p_n$. We are going to show that there are three cases depending on \vec{p}_n :

1. If $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = \Theta(\sum_{i=2}^n p_i^2)$, then we can show that the asymptotic satisfiability threshold is at $m = \Theta(q_{\max}^{-1})$, where $q_{\max} = \Theta((p_1 \cdot p_2) / (1 - \sum_{i=1}^n p_i^2))$ is the maximum clause probability. We can also show that this threshold is coarse. The coarseness stems from the emergence of an unsatisfiable sub-formula of size 4, which contains only the two most probable variables.
2. If $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = o(\sum_{i=2}^n p_i^2)$, then the asymptotic threshold is at $m = \Theta\left(\left(1 - \sum_{i=1}^n p_i^2\right) / \left(p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)\right)$ and it is again coarse. This time the coarseness stems from the emergence of an unsatisfiable sub-formula with 3 variables and 4 clauses.
3. If $p_1^2 = o(\sum_{i=1}^n p_i^2)$, then there is a sharp threshold at exactly $m = 1 / \left(\sum_{i=1}^n p_i^2\right)$.

Note that these three cases give us a complete dichotomy of coarseness and sharpness for the satisfiability threshold of non-uniform random 2-SAT. This result generalizes the seminal works by Chvatal and Reed [11] and by Goerdt [24] to arbitrary variable probability distributions and includes their findings as a special case (c.f. Section 6). We summarize our findings in the following theorem.

Theorem 1.1. *Let $\mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ be the non-uniform random 2-SAT model with n variables, m clauses, and an ensemble of probability distributions $(\vec{p}_x)_{x \in \mathbb{N}}$. Let $\vec{p}_n = (p_1, p_2, \dots, p_n)$ be the n -th distribution from the ensemble. W.l.o.g. let $p_1 \geq p_2 \geq \dots \geq p_n$. If $p_1^2 = o(\sum_{i=1}^n p_i^2)$, then $\mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ has a sharp satisfiability threshold at $m = (\sum_{i=1}^n p_i^2)^{-1}$. Otherwise, $\mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ has a coarse satisfiability threshold at $m = \Theta\left(\left(1 - \sum_{i=1}^n p_i^2\right) / \left(p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)\right)$.*

1.4 Techniques

For the sharp threshold result, we only show the conditions on sharpness. These also imply the existence of an asymptotic threshold. For the coarse threshold results, however, we first have to show the existence of an asymptotic threshold at some number of clauses $m^*(n)$. Then, we have to show that for some range of constants $\varepsilon \in [\varepsilon_1, \varepsilon_2]$ the probability to generate a satisfiable instance at $\varepsilon \cdot m^*(n)$ is a constant bounded away from zero and one.

We extend and generalize the proof ideas of Chvatal and Reed [11]. In order to show a lower bound on the threshold, we investigate the existence of bicycles. Bicycles were introduced by Chvatal and Reed. They are sub-formulas which appear in every unsatisfiable formula. We can show with a first moment argument, that these do not appear below a certain number of clauses, thus making formulas satisfiable.

In order to show an upper bound on the threshold, we investigate the existence of snakes. Snakes are unsatisfiable sub-formulas and have also been introduced by Chvatal and Reed. We can show with a second-moment argument that snakes of certain sizes do appear above a certain number of clauses, thus making formulas unsatisfiable. However, we need to be careful and distinguish more possibilities of partially mapping snakes onto each other than in the uniform case. Unfortunately, this method does not work if $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = \Theta(\sum_{i=2}^n p_i^2)$. In that case we lower-bound the probability that an unsatisfiable sub-formula containing only the two most-probable variables exists. This can be done with a simple inclusion-exclusion argument and the resulting lemma also work for $k \geq 3$.

2 Preliminaries

We analyze non-uniform random k -SAT on n variables and m clauses. We denote by X_1, \dots, X_n the Boolean variables. A clause is a disjunction of k literals $\ell_1 \vee \dots \vee \ell_k$, where each literal assumes a (possibly negated) variable. For a literal ℓ_i let $|\ell_i|$ denote the variable of the literal. A formula Φ in conjunctive normal form is a conjunction of clauses $c_1 \wedge \dots \wedge c_m$. We conveniently interpret a clause c both as a Boolean formula and as a set of literals. We say that Φ is satisfiable if there exists an assignment of variables X_1, \dots, X_n such that the formula evaluates to 1. Now let $(\vec{p}_n)_{n \in \mathbb{N}}$ be an ensemble of probability distributions, where $\vec{p}_n = (p_{n,1}, p_{n,2}, \dots, p_{n,n})$ is a probability distribution over n variables with $\Pr(X = X_i) = p_{n,i} =: p_n(X_i)$.

Definition 2.1 (Clause-Drawing Non-Uniform Random k -SAT). *Let m, n, k be given, and consider any ensemble of probability distributions $(\vec{p}_n)_{n \in \mathbb{N}}$, where $\vec{p}_n = (p_{n,1}, p_{n,2}, \dots, p_{n,n})$ is a probability distribution over n variables with $\sum_{i=1}^n p_{n,i} = 1$. The clause-drawing non-uniform random k -SAT (non-uniform random k -SAT) model $\mathcal{D}(n, k, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ constructs a random SAT formula Φ by sampling m clauses independently at random. Each clause is sampled as follows:*

1. *Select k variables independently at random from the distribution \vec{p}_n . Repeat until no variables coincide.*
2. *Negate each of the k variables independently at random with probability $1/2$.*

For the sake of simplicity and since we will always only consider one distribution from the ensemble, we will omit the index n throughout the paper, e.g. the probability distribution \vec{p}_n will be denoted as (p_1, p_2, \dots, p_n) . W.l.o.g. we will assume $p_1 \geq p_2 \geq \dots \geq p_n$.

The clause-drawing non-uniform random k -SAT model is equivalent to drawing each clause independently at random from the set of all k -clauses which contain no variable more than once. The probability to draw a clause c over n variables is then

$$q_c := \frac{\prod_{\ell \in c} p(|\ell|)}{2^k \sum_{J \in \mathcal{P}_k(\{1,2,\dots,n\})} \prod_{j \in J} p_j}, \quad (2.1)$$

where $\mathcal{P}_k(\cdot)$ denotes the set of cardinality- k elements of the power set. The factor 2^k in the denominator comes from the different possibilities to negate variables. Note that $k! \sum_{J \in \mathcal{P}_k(\{1,2,\dots,n\})} \prod_{j \in J} p_{n,j}$ is the probability of choosing a k -clause that contains no variable more than once. We can now write

$$q_c = C \frac{k!}{2^k} \prod_{X \in S} p_n(X), \quad (2.2)$$

where we define $C := 1 / \left(k! \cdot \sum_{J \in \mathcal{P}_k(\{1,2,\dots,n\})} \prod_{j \in J} p_{n,j} \right)$. For $k = 2$ it holds that $C = 1 / (1 - (\sum_{i=1}^n p_i^2))$. Hiding this factor in C makes clause probabilities easier to handle. Throughout the paper we let q_{\max} denote the maximum clause probability as defined in Equation (2.2). In Section 3 and Section 4 we will assume $q_{\max} = o(1)$. The case $q_{\max} = \Theta(1)$ will be handled in Section 5. Note that this case can only happen for $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = \Theta(\sum_{i=2}^n p_i^2)$.

3 Bi-Cycles and a Lower Bound on the Satisfiability Threshold

Chvatal and Reed [11] define the following sub-structure of 2-SAT formulas and show that every unsatisfiable 2-CNF contains this substructure.

Definition 3.1 (bi-cycle). *We define a bicycle of length t to be a sequence of $t + 1$ clauses of the form*

$$(u, w_1), (\bar{w}_1, w_2), \dots, (\bar{w}_{t-1}, w_t), (\bar{w}_t, v),$$

where w_1, \dots, w_t are literals of distinct variables and $u, v \in \{w_1, \dots, w_t, \bar{w}_1, \dots, \bar{w}_t\}$.

To lower-bound the probability for a random 2-CNF to be satisfiable it is therefore sufficient to upper-bound the probability that such a formula contains a bicycle. This is done in the following two lemmas. Their proofs are oriented along the lines of the proof of Theorem 3 from [11].

Lemma 3.1. *Consider a non-uniform random 2-SAT formula Φ with $p_1^2 = o(\sum_{i=1}^n p_i^2)$. Then, Φ is satisfiable with probability at least $1 - o(1)$ for a number of clauses $m < (1 - \varepsilon) (\sum_{i=1}^n p_i^2)^{-1}$, where $\varepsilon > 0$ is a constant.*

Proof. To show this result, we show that the expected number of bicycles is $o(1)$ for the setting we consider. The result then follows by Markov's inequality.

First, we fix a set $S \subseteq [n]$ of variables to appear in a bicycle with $|S| = t \geq 2$. The probability that a *specific* bicycle B with these variables appears in Φ is

$$\Pr(B \text{ in } \Phi) = \underbrace{\binom{m}{t+1}}_{\text{positions of } B \text{ in } \Phi} (t+1)! \cdot \Pr(u \vee w_1) \cdot \Pr(\bar{w}_t \vee v) \prod_{h=1}^{t-1} \Pr(\bar{w}_h \vee w_{h+1}).$$

For literals w_i over variables x_i it holds that

$$\Pr(w_j \vee w_i) = \frac{C}{2} p_i \cdot p_j,$$

where $1 \leq C = (1 - \sum_{i=1}^n p_i^2)^{-1} = 1 + o(1)$, since $\sum_{i=1}^n p_i^2 = o(1)$ due to the requirement $p_1^2 = o(\sum_{i=1}^n p_i^2)$. There are at most $t!$ possibilities to arrange the t variables in a bicycle and 2^t possibilities to choose literals from the t variables. For the probability that *any* bicycle with the variables from S appears in Φ it now holds that

$$\Pr(S\text{-bicycle in } \Phi) \leq m^{t+1} \cdot t! \cdot 2^t \cdot \left(\frac{C}{2}\right)^{t+1} \cdot \prod_{i \in S} p_i^2 \left(2 \cdot \sum_{i \in S} p_i\right)^2$$

where the last factor accounts for the possibilities to choose u and v . It now holds that

$$\begin{aligned} \Pr(\Phi \text{ contains a bicycle}) &\leq \sum_{t=2}^n \sum_{S \subseteq \mathcal{P}_t(V)} m^{t+1} \cdot t! \cdot 2^t \cdot \left(\frac{C}{2}\right)^{t+1} 2^2 \cdot \prod_{i \in S} p_i^2 \left(\sum_{i \in S} p_i\right)^2 \\ &\leq 2 \cdot \sum_{t=2}^n (C \cdot m)^{t+1} \cdot t! \cdot t^2 \cdot p_1^2 \cdot \sum_{S \subseteq \mathcal{P}_t(V)} \prod_{i \in S} p_i^2 \\ &\leq 2 \cdot \sum_{t=2}^n (C \cdot m)^{t+1} \cdot t^2 \cdot p_1^2 \cdot \left(\sum_{i \in S} p_i^2\right)^t \\ &= o\left(2 \cdot \sum_{t=2}^n \left(C \cdot m \left(\sum_{i \in S} p_i^2\right)\right)^{t+1} \cdot t^2\right), \end{aligned}$$

where we used $\sum_{i \in S} p_i \leq t \cdot p_1$ in the second, $\sum_{S \subseteq \mathcal{P}_t(V)} \prod_{i \in S} p_i^2 \leq \frac{1}{t!} \cdot (\sum_{i \in S} p_i^2)^t$ in the third line, and the requirement $p_1^2 = o(\sum_{i=1}^n p_i^2)$ in the fourth line. It is obvious that this probability is $o(1)$ as soon as the sum becomes a constant. This holds for $m < (1 - \varepsilon) (\sum_{i=1}^n p_i^2)^{-1} < (C \cdot \sum_{i=1}^n p_i^2)^{-1}$, where $\varepsilon > 0$ is a constant. \square

It has to be noted that in the former lemma we ignored the factor C in our bound. We can do this, since for $p_1^2 = o(\sum_{i=1}^n p_i^2)$ it always is $1 + o(1)$ and does not make a difference for sharpness due to our definition. In the case of $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$, we can show the following result with a similar proof, but now we have to take C into account, since it might become super-constant.

Lemma 3.2. *Consider a non-uniform random 2-SAT formula Φ with $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $q_{\max} = o(1)$. Then, Φ is satisfiable with probability at least $1 - o(1)$ for a number of clauses $m = o\left(\left(C \cdot p_1 \cdot (\sum_{i=2}^n p_i^2)^{1/2}\right)^{-1}\right)$. Also, there is a constant $\varepsilon \in (0, 1)$ such that Φ is satisfiable with a positive constant probability for a number of clauses $m \leq (1 - \varepsilon) \left(C \cdot p_1 \cdot (\sum_{i=2}^n p_i^2)^{1/2}\right)^{-1}$.*

Proof. As in the proof of Lemma 3.1 it holds that

$$\begin{aligned} \Pr(\Phi \text{ unsat}) &\leq \Pr(\Phi \text{ contains a bicycle}) \\ &\leq \sum_{t=2}^n \sum_{S \subseteq \mathcal{P}_t(V)} m^{t+1} \cdot t! \cdot 2^t \cdot \left(\frac{C}{2}\right)^{t+1} 2^2 \cdot \prod_{i \in S} p_i^2 \left(\sum_{i \in S} p_i\right)^2 \\ &\leq 2 \cdot \sum_{t=2}^n (C \cdot m)^{t+1} \cdot t! \cdot \sum_{S \subseteq \mathcal{P}_t(V)} \left(\prod_{i \in S} p_i^2\right) \cdot \left(\sum_{i \in S} p_i\right)^2. \end{aligned} \quad (3.1)$$

We can now do a more detailed analysis of the term $\sum_{S \subseteq \mathcal{P}_t(V)} \left(\prod_{i \in S} p_i^2\right) \cdot \left(\sum_{i \in S} p_i\right)^2$ as follows

$$\begin{aligned} &\sum_{S \subseteq \mathcal{P}_t(V)} \left(\prod_{i \in S} p_i^2\right) \cdot \left(\sum_{i \in S} p_i\right)^2 \\ &\leq p_1^2 \cdot t^2 \cdot p_1^2 \cdot \frac{1}{(t-1)!} \cdot \left(\sum_{i=2}^n p_i^2\right)^{t-1} + t^2 \cdot p_2^2 \cdot \frac{1}{t!} \cdot \left(\sum_{i=2}^n p_i^2\right)^t \\ &= \mathcal{O}\left(t^3 \cdot p_1^4 \cdot \frac{1}{t!} \cdot \left(\sum_{i=2}^n p_i^2\right)^{t-1}\right) \end{aligned} \quad (3.2)$$

where the second line is just a case distinction between the terms with $p_1 \in S$ and $p_1 \notin S$ and the last line follows due to $p_2 \leq p_1$ and the requirement $\sum_{i=2}^n p_i^2 \leq \sum_{i=1}^n p_i^2 = \mathcal{O}(p_1^2)$. It holds that $p_1^4 \cdot (\sum_{i=2}^n p_i^2)^{t-1} = \mathcal{O}\left(\left(p_1 \cdot (\sum_{i=2}^n p_i^2)^{1/2}\right)^{t+1}\right)$ for $t \geq 3$. For $t = 2$ we can actually show that

$$\begin{aligned}
& \sum_{S \subseteq \mathcal{P}_x(V)} \Pr(S\text{-bicycle in } F) \\
& \leq (C \cdot m)^3 \cdot \sum_{i,j \in V} p_i^3 \cdot p_j^3 \\
& \leq (C \cdot m)^3 \cdot p_1^3 \cdot \left(\sum_{i=2}^n p_i^3\right) + (C \cdot m)^3 \left(\sum_{i=2}^n p_i^3\right)^2 \\
& \leq (C \cdot m)^3 \cdot p_1^3 \cdot \left(\sum_{i=2}^n p_i^2\right)^{3/2} + (C \cdot m)^3 \cdot \left(\sum_{i=2}^n p_i^2\right)^3 \\
& = \mathcal{O}\left((C \cdot m)^3 \cdot \left(p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{t+1}\right) \tag{3.3}
\end{aligned}$$

where the first line holds since each of the three 2-clauses in the bicycle must contain both variables, the second line is again a case distinction, the third line follows due to the monotonicity of vector norms, and the fourth line follows due to $p_2 \leq p_1$ and $\sum_{i=2}^n p_i^2 \leq \sum_{i=1}^n p_i^2 = \mathcal{O}(p_1^2)$.

We can now plug Equation (3.2) and Equation (3.3) into Equation (3.1) to get

$$\Pr(\Phi \text{ unsat}) \leq 2 \cdot K \cdot \sum_{t=2}^n \left(C \cdot m \cdot \left(p_1 \cdot \left(\sum_{i=2}^n p_i^2 \right)^{1/2} \right) \right)^{t+1} t^3.$$

for some constant K that is only determined by the probability vector. That means, for $m = o\left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}$ the expression is $o(1)$ and for $m \leq (1 - \varepsilon) \cdot \left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}$ for some sufficiently large constant $\varepsilon \in (0, 1)$, the expression is a constant smaller than 1 as desired. \square

Note that this lemma captures both cases for $p_1^2 = \Theta\left(\sum_{i=1}^n p_i^2\right)$. If also $p_2^2 = \Theta\left(\sum_{i=2}^n p_i^2\right)$, then $\left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1} = \Theta(q_{\max}^{-1})$ is the asymptotic threshold as we stated in the introduction. The case $q_{\max} = \Theta(1)$ has to be excluded, since for that case the asymptotic threshold is a constant. The above lemma might then give us a value so small that the ranges where we can lower- and upper-bound satisfiability to constants away from zero resp. one do not overlap. Thus, this case is handled separately in Section 5.

4 Snakes and an Upper Bound on the Satisfiability Threshold

The two lemmas from the previous section provided a lower bound on the satisfiability threshold for non-uniform random 2-SAT. By using the second moment method, we can also derive an upper bound on the threshold. Again, this proof is inspired by Chvatal and Reed [11, Theorem 4], who provide us with the following definition.

Definition 4.1 (snake). *A snake of size t is a sequence of literals $w_1, w_2, \dots, w_{2t-1}$ over distinct variables. Each snake A is associated with a set F_A of $2t$ clauses (\bar{w}_i, w_{i+1}) , $0 \leq i \leq 2t - 1$, such that $w_0 = w_{2t} = \bar{w}_t$.*

We will also call the variable $|w_t|$ of a snake its *central* variable. Note that the set of clauses F_A defined by a snake A is unsatisfiable. Also, the snakes $(w_1, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_s)$, $(\bar{w}_{t-1}, \bar{w}_{t-2}, \dots, \bar{w}_1, w_t, w_{t+1}, \dots, w_s)$, $(w_1, \dots, w_{t-1}, w_t, \bar{w}_s, \bar{w}_{s-1}, \dots, \bar{w}_{t+1})$ and $(\bar{w}_{t-1}, \bar{w}_{t-2}, \dots, \bar{w}_1, w_t \bar{w}_s, \bar{w}_{s-1}, \dots, \bar{w}_{t+1})$ create the same set of formulas.

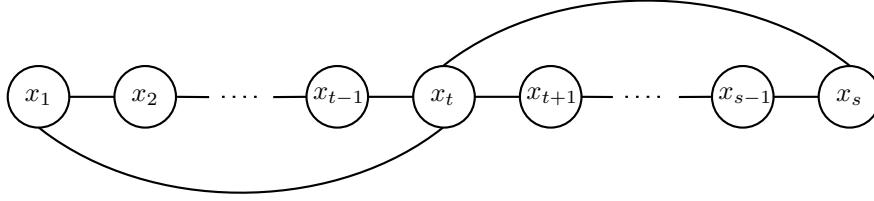


Figure 1: Variable-variable-incidence graph of a snake w_1, w_2, \dots, w_s where $|w_i| = x_i$ (the variable of the literal w_i) for $1 \leq i \leq s = 2t - 1$.

The *variable-variable incidence graph (VIG)* for a formula Φ is a simple graph $G_\Phi = (V_\Phi, E_\Phi)$ with V_Φ consisting of all variables appearing in Φ and two variables being connected by an edge if they appear together in at least one clause of Φ . An example for a snake's VIG can be seen in Figure 1. This representation will come in handy later in the proof of Lemma 4.6.

In order to show our upper bounds, we will prove that snakes of a certain length t appear with sufficiently high probability in a random formula $\Phi \sim \mathcal{D}(n, k, (\vec{p}_x)_{x \in \mathbb{N}}, m)$. To this end we utilize the second moment method: If $X \geq 0$ is a random variable with finite variance, then

$$\Pr(X > 0) \geq \frac{\mathbb{E}[X]^2}{\mathbb{E}[X^2]}.$$

We define the following indicator variables for each snake A of size t

$$X_A = \begin{cases} 1 & \text{if } F_A \text{ appears exactly once in } \Phi \\ 0 & \text{otherwise} \end{cases}$$

and their sum $X_t = \sum_A X_A$. For carefully chosen t we will show $\mathbb{E}[X_t^2] = \mathcal{O}(\mathbb{E}[X_t]^2)$ to show a coarse and $\mathbb{E}[X_t^2] = (1 + o(1)) \cdot (\mathbb{E}[X_t]^2)$ to show a sharp threshold. This implies a constant resp. $1 - o(1)$ probability to be unsatisfiable due to the second moment method. In the case of $p_1^2 = o(\sum_{i=1}^n p_i^2)$, we will choose $t = \Theta(\log^2 f(n))$, where we define $f(n) = (\sum_{i=1}^n p_i^2) / p_1^2$. For $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = o(\sum_{i=2}^n p_i^2)$ we choose $t = 2$. We only want to use the method for these two cases. The third case with $p_1 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2 = \Theta(\sum_{i=2}^n p_i^2)$ will be handled with the more general Lemma 4.7.

Now, if we want to use the second moment method, we first have to ensure that the expected number of snakes of a certain size is large enough. The following lemma provides a lower bound on this expected number.

Lemma 4.1. *Let X_t be the number of snakes of size $s + 1 = 2t$ whose associated formulas appear exactly once in a non-uniform random 2-SAT formula. Then it holds that*

$$\mathbb{E}[X_t] \geq \frac{1}{2} (m - 2t)^{2t} \cdot C^{2t} \cdot e^{-(m-2t) \frac{2t \cdot q_{\max}}{1-2t \cdot q_{\max}}} \cdot \left(\sum_{i=1}^n p_i^4 \right) \cdot \left(\sum_{i=2}^n p_i^2 - (2t-2) \cdot p_2^2 \right)^{2t-2}.$$

Proof. It holds that

$$\begin{aligned} & \mathbb{E}[X_t] \tag{4.1} \\ &= \sum_{\text{snake } A=(w_1, \dots, w_{2t-1})} \binom{m}{2t} \cdot (2t)! \cdot \prod_{i=0}^{2t-1} \Pr((\bar{w}_i, w_{i+1})) \cdot \left(1 - \sum_{c \in F_A} P(c) \right)^{m-2t} \\ &\geq (m-2t)^{2t} \left(1 - \sum_{c \in F_A} P(c) \right)^{m-2t} \left(\frac{C}{2} \right)^{2t} \sum_{\text{snake } A=(w_1, \dots, w_s)} p(|w_i|)^4 \cdot \prod_{\substack{i=1 \\ i \neq t}}^{2t-1} p(|w_i|)^2 \\ &\geq (m-2t)^{2t} \left(1 - \sum_{c \in F_A} P(c) \right)^{m-2t} \left(\frac{C}{2} \right)^{2t} 2^{2t-1} \sum_{j=1}^n \left(p_j^4 \cdot (2t-2)! \cdot \sum_{S \subseteq [n] \setminus \{j\}} \prod_{s \in S} p_s^2 \right). \tag{4.2} \end{aligned}$$

First, notice that

$$\sum_{S \subseteq [n] \setminus \{j\}} \prod_{s \in S} p_s^2 \geq \frac{1}{(2t-2)!} \left(\sum_{i=2}^n p_i^2 - (2t-2) \cdot p_2^2 \right)^{2t-2}. \tag{4.3}$$

It now holds that

$$\left(1 - \sum_{c \in F_A} P(c)\right)^{m-2t} \geq (1 - 2t \cdot q_{\max})^{m-2t} > \exp\left(- (m-2t) \frac{2t \cdot q_{\max}}{1 - 2t \cdot q_{\max}}\right), \quad (4.4)$$

where we used $(1-x) > e^{-\frac{x}{1-x}}$ for $x \in [0, 1)$. Plugging Equation (4.3) and Equation (4.4) into Equation (4.2) we get the result as desired. \square

In order to use the second moment method we have to show that this expected value is at least a constant if we want to show a coarse threshold and asymptotically bigger than a constant if we want to show a sharp threshold. Hence, the following lemmas give lower bounds on $\mathbb{E}[X_t]$ for the first two cases and the respective ranges of t we consider.

Lemma 4.2. *Let X_t be the number of snakes of size t that appear exactly once in a non-uniform random 2-SAT formula with $p_1^2 = o(\sum_{i=1}^n p_i^2)$ and $m = (1 + \varepsilon) (\sum_{i=1}^n p_i^2)^{-1}$ for some $\varepsilon > 0$. Then it holds that*

$$\mathbb{E}[X_t] \geq (1 - o(1)) \cdot m^{2t} \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=1}^n p_i^2\right)^{2t-2} = \omega(1)$$

if $t = o(\sqrt{f(n)}) \cap \omega(\log f(n))$, where $f(n) = (\sum_{i=1}^n p_i^2) / p_1^2$.

Proof. It holds that

$$\mathbb{E}[X_t] \geq \frac{1}{2} (m-2t)^{2t} \cdot C^{2t} \cdot e^{-(m-2t) \frac{2t \cdot q_{\max}}{1-2t \cdot q_{\max}}} \cdot \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=2}^n p_i^2 - (2t-2) \cdot p_1^2\right)^{2t-2}.$$

Furthermore,

$$\begin{aligned} \left(\sum_{i=2}^n p_i^2 - (2t-2)p_1^2\right)^{2t-2} &\geq \left(\sum_{i=1}^n p_i^2 - (2t-3) \cdot p_1^2\right)^{2t-2} \\ &= \left(\sum_{i=1}^n p_i^2\right)^{2t-2} \cdot \left(1 - \frac{(2t-3) \cdot p_1^2}{\sum_{i=1}^n p_i^2}\right)^{2t-2} \\ &\geq \left(\sum_{i=1}^n p_i^2\right)^{2t-2} \cdot \exp\left(-\frac{(2t-2) \cdot (2t-3) / f(n)}{1 - (2t-3) / f(n)}\right) \\ &= \left(\sum_{i=1}^n p_i^2\right)^{2t-2} \cdot (1 - o(1)), \end{aligned}$$

where we used $t = o(\sqrt{f(n)})$ in the last line. Equivalently,

$$\begin{aligned} (m-2t)^{2t} &\geq m^{2t} \cdot \exp\left(-\frac{4t^2/m}{1-2t/m}\right) \\ &= m^{2t} \cdot (1 - o(1)), \end{aligned}$$

which holds since $t^2 = o(f(n))$ and $f(n) \cdot p_1^2 = \sum_{i=1}^n p_i^2 \leq \sum_{i=1}^n p_1 \cdot p_i = p_1$ implies $f(n) \leq p_1^{-1}$ and thus $m = 1 / (\sum_{i=1}^n p_i^2) = 1 / (f(n) \cdot p_1^2) \geq f(n)$. Since we know that $C = \frac{1}{1 - \sum_{i=1}^n p_i^2}$, this also implies $1 \leq C \leq \frac{1}{1-1/f(n)} = \mathcal{O}(1)$. We also know that

$$\exp\left(- (m-2t) \frac{2t \cdot q_{\max}}{1 - 2t \cdot q_{\max}}\right) = 1 - o(1),$$

as $q_{\max} = \mathcal{O}(C \cdot p_1^2)$ implies $m \cdot t \cdot q_{\max} = \mathcal{O}\left(\frac{t \cdot p_1^2}{\sum_{i=1}^n p_i^2}\right) = \mathcal{O}\left(\frac{t}{f(n)}\right) = o(1)$. The expected value now simplifies to

$$\mathbb{E}[X_t] = (1 - o(1)) \cdot m^{2t} \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=1}^n p_i^2\right)^{2t-2}.$$

It holds that

$$m^2 \cdot \left(\sum_{i=1}^n p_i^4 \right) \geq m^2 \cdot p_1^4 = \frac{(1+\varepsilon)^2 \cdot p_1^4}{\left(\sum_{i=1}^n p_i^2 \right)^2} = \frac{(1+\varepsilon)^2}{f(n)^2},$$

where we used $m = (1+\varepsilon) \left(\sum_{i=1}^n p_i^2 \right)^{-1}$. With the same fact it holds that

$$\left(m \cdot \sum_{i=1}^n p_i^2 \right)^{2t-2} = (1+\varepsilon)^{2t-2}.$$

Since we know that $t = \omega(\log f(n))$, it holds that

$$\mathbb{E}[X_t] \geq (1-o(1)) \cdot \frac{(1+\varepsilon)^{2t}}{f(n)^2} = \omega(1)$$

as desired. \square

Lemma 4.3. *Let X_t be the number of snakes of size t that appear exactly once in a non-uniform random 2-SAT formula with $p_1^2 = \Theta\left(\sum_{i=1}^n p_i^2\right)$ and $p_2^2 = o\left(\sum_{i=2}^n p_i^2\right)$. For $t = 2$ and $m = \Omega\left(\left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right) \cap o\left(q_{\max}^{-1}\right)$ it holds that*

$$\mathbb{E}[X_2] \geq (1-o(1)) \cdot m^4 \cdot C^4 \cdot p_1^4 \cdot \left(\sum_{i=1}^n p_i^2 \right)^2.$$

Furthermore,

$$\mathbb{E}[X_2] = \begin{cases} \Omega(1) & , m = \Theta\left(\left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right) \text{ and} \\ \omega(1) & , m = \omega\left(\left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right) \cap o\left(\left(q_{\max}\right)^{-1}\right). \end{cases}$$

Proof. First, note that the range of m in the second case is not empty, since

$$q_{\max}^{-1} = \Omega\left(\frac{1}{C \cdot p_1 \cdot p_2}\right) = \omega\left(\frac{1}{C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}}\right)$$

due to $p_2^2 = o\left(\sum_{i=2}^n p_i^2\right)$. With $t = 2$ it holds that

$$\mathbb{E}[X_2] \geq \frac{1}{2}(m-4)^4 \cdot C^4 \cdot e^{-(m-4) \frac{4 \cdot q_{\max}}{1-4 \cdot q_{\max}}} \cdot \left(\sum_{i=1}^n p_i^4 \right) \cdot \left(\sum_{i=2}^n p_i^2 - 2 \cdot p_2^2 \right)^2$$

due to Lemma 4.1 We now get

$$\left(\sum_{i=2}^n p_i^2 - 2 \cdot p_2^2 \right)^2 \geq \left(\sum_{i=2}^n p_i^2 \right)^2 \cdot \left(1 - \frac{2 \cdot p_2^2}{\sum_{i=2}^n p_i^2} \right)^2 = \left(\sum_{i=2}^n p_i^2 \right)^2 \cdot (1-o(1)),$$

where we used $p_2^2 = o\left(\sum_{i=2}^n p_i^2\right)$. Equivalently,

$$(m-4)^4 \geq m^4 \cdot \left(1 - \frac{4}{m} \right)^4 = m^4 \cdot (1-o(1)),$$

which holds since $m = \Omega\left(C \cdot p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1} = \omega(1)$. First, we can see that

$$\sum_{i=2}^n p_i^2 \leq p_2 \cdot \sum_{i=2}^n p_i = o\left(\left(\sum_{i=2}^n p_i^2\right)^{1/2} \cdot \left(\sum_{i=2}^n p_i\right)\right) = o\left(\left(\sum_{i=2}^n p_i\right)^2\right),$$

since $p_2^2 = o(\sum_{i=2}^n p_i^2)$ and since $\sum_{i=2}^n p_i^2 \leq (\sum_{i=2}^n p_i)^2$. Now we distinguish two cases. Either $p_1 = 1 - \varepsilon$ for some constant $\varepsilon > 0$ or $p_1 = 1 - 1/g(n)$ for some $g(n) = \omega(1)$. In the first case, $C = \mathcal{O}(1)$. This implies

$$m^{-1} = \mathcal{O}\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right) = o\left(C \cdot \left(\sum_{i=2}^n p_i\right)\right) = o(1),$$

since $p_1 \leq 1$ and $\sum_{i=2}^n p_i = \mathcal{O}(1)$. In the second case, $C = \mathcal{O}(g(n))$, but also $\sum_{i=2}^n p_i \leq 1/g(n)$. Thus,

$$m^{-1} = \mathcal{O}\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right) = o\left(C \cdot \left(\sum_{i=2}^n p_i\right)\right) = o(1).$$

This gives us $m = \omega(1)$ and it implies $q_{\max} = o(1)$. We know that

$$\exp\left(- (m-4) \frac{4 \cdot q_{\max}}{1 - 4 \cdot q_{\max}}\right) = 1 - o(1),$$

as $m = o(q_{\max}^{-1})$. The expected value now simplifies to

$$\mathbb{E}[X_2] = (1 - o(1)) \cdot m^4 \cdot C^4 \cdot \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=1}^n p_i^2\right)^2 \geq (1 - o(1)) \cdot m^4 \cdot C^4 \cdot p_1^4 \cdot \left(\sum_{i=1}^n p_i^2\right)^2,$$

since $(\sum_{i=1}^n p_i^4) \geq p_1^4$. It holds that $\mathbb{E}[X_2] = \Omega(1)$ for $m = \Theta\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right)$ and $\mathbb{E}[X_2] = \omega(1)$ for $m = \omega\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right)$ as desired. \square

Now we are ready to prove an upper bound on the non-uniform random 2-SAT threshold. To get to know the proof technique, we start with the much simpler case $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = o(\sum_{i=2}^n p_i^2)$. The proof contains a small case distinction depending on how the shared clauses of two snakes A and B influence $\Pr(X_A \wedge X_B)$. The next lemma establishes that there is a regime of m where random formulas are unsatisfiable with a positive constant probability.

Lemma 4.4. *Consider a non-uniform random 2-SAT formula Φ with $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = o(\sum_{i=2}^n p_i^2)$. Then Φ is unsatisfiable with positive constant probability for $m = \Theta\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right)$.*

Proof. First, we want to show that for $m = \Theta\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right)$, F_A for a snake A of size $|F_A| = 4$ appears in Φ with constant probability. Since Lemma 4.3 gives us a lower bound on $\mathbb{E}[X_2]$, we only need to consider $\mathbb{E}[X_2^2]$ now. We use the same approach as Chvatal and Reed [11] and split $\mathbb{E}[X_2^2]$ into two parts as follows

$$\mathbb{E}[X_2^2] = \sum_A \sum_B \Pr(X_A \wedge X_B) = \sum_A \left(\sum_{B: B \not\sim A} \Pr(X_A \wedge X_B) + \sum_{B: B \sim A} \Pr(X_A \wedge X_B) \right),$$

where $B \sim A$ denotes $F_A \cap F_B \neq \emptyset$. We will show that the part for $B \sim A$ is at most $(1 + o(1)) \cdot \mathbb{E}[X_2]^2$ and that the other part is $\mathcal{O}(\mathbb{E}[X_2]^2)$.

First let us consider the part for $B \sim A$. It holds that

$$\Pr(X_A \wedge X_B) = \binom{m}{8} \cdot 8! \cdot \left(\prod_{c \in F_A} \Pr(c)\right) \cdot \left(\prod_{c \in F_B} \Pr(c)\right) \cdot \left(1 - \sum_{c \in F_A \cup F_B} \Pr(c)\right)^{m-8},$$

while

$$\Pr(X_A) = \binom{m}{4} \cdot 4! \cdot \left(\prod_{c \in F_A} \Pr(c)\right) \cdot \left(1 - \sum_{c \in F_A} \Pr(c)\right)^{m-4}. \quad (4.5)$$

This readily implies

$$\begin{aligned} \Pr(X_A \wedge X_B) &\leq \Pr(X_A) \cdot \Pr(X_B) \frac{(1 - \sum_{c \in F_A \cup F_B} \Pr(c))^{m-8}}{(1 - \sum_{c \in F_A} \Pr(c))^{m-4} (1 - \sum_{c \in F_B} \Pr(c))^{m-4}} \\ &\leq (1 + o(1)) \cdot \Pr(X_A) \cdot \Pr(X_B), \end{aligned}$$

since $\binom{m}{8} \cdot 8! \leq ((\binom{m}{4} \cdot 4!))^2$ and since

$$\left(1 - \sum_{c \in F_A} \Pr(c)\right)^{m-4} \geq (1 - 4 \cdot q_{\max})^{m-4} > \exp\left(-\frac{4 \cdot (m-4) \cdot q_{\max}}{1 - 4 \cdot q_{\max}}\right) = 1 - o(1)$$

for any snake A with 4 clauses, since $m \cdot q_{\max} = \mathcal{O}\left(\frac{p_1 \cdot p_2}{p_1 \cdot (\sum_{i=2}^n p_i^2)^{1/2}}\right) = o(1)$. This establishes

$$\sum_A \sum_{B: B \approx A} \Pr(X_A \wedge X_B) \leq (1 + o(1)) \sum_A \sum_{B: B \approx A} \Pr(X_A) \Pr(X_B) \leq (1 + o(1)) \mathbb{E}[X_2]^2. \quad (4.6)$$

Now we turn to the case that $B \sim A$. We want to show that this second sum is $\mathcal{O}(\mathbb{E}[X_2]^2)$. Let $l = |F_A \cap F_B|$. The first and simplest case is $F_A = F_B$. This obviously happens if $A = B$, but also for three other snakes. So it holds that

$$\sum_A \sum_{B: |F_A \cap F_B|=4} \Pr(X_A \wedge X_B) = 4 \cdot \mathbb{E}[X_2] = \mathcal{O}(\mathbb{E}[X_2]^2), \quad (4.7)$$

since $\Pr(X_A \wedge X_B) = \Pr(X_A)$ and $\mathbb{E}[X_2] = \Omega(1)$.

It now holds that

$$\begin{aligned} &\sum_A \sum_{B: |F_A \cap F_B|=l} \Pr(X_A \wedge X_B) \\ &\leq \binom{m}{8-l} \cdot (8-l)! \cdot \left(1 - \sum_{c \in F_A \cup F_B} \Pr(c)\right)^{m-8+l} \cdot 2^3 \cdot 2! \left(\frac{C}{2}\right)^4 \\ &\cdot \left(\sum_{i=1}^n p_i^4 \cdot \sum_{\substack{S \subseteq [n] \setminus \{i\}: \\ |S|=2}} \prod_{s \in S} p_s^2\right) \cdot \sum_{B: |F_A \cap F_B|=l} \prod_{c \in F_B \setminus F_A} \Pr(c) \end{aligned} \quad (4.8)$$

where we accounted for the l possible positions of clauses from $F_A \cup F_B$ in Φ , for the $2^3 \cdot 2!$ possibilities to create a snake A from chosen variables if the central variable is determined already, and for the ways to choose those variables. Now we want to bound $\sum_{i=1}^n p_i^4 \cdot \sum_{\substack{S \subseteq [n] \setminus \{i\}: \\ |S|=2}} \prod_{s \in S} p_s^2$.

In order to do so we distinguish between the cases that p_1 appears in the snake as the central variable, a non-central variable or not at all to show the following

$$\begin{aligned} &\sum_{i=1}^n \left(p_i^4 \cdot \sum_{\substack{S \subseteq [n] \setminus \{i\}: \\ |S|=2}} \prod_{s \in S} p_s^2 \right) \\ &\leq p_1^4 \cdot \left(\sum_{i=2}^n p_i^2\right)^2 + \left(\sum_{i=2}^n p_i^4\right) \cdot p_1^2 \cdot \left(\sum_{i=2}^n p_i^2\right) + \left(\sum_{i=2}^n p_i^4\right) \cdot \left(\sum_{i=2}^n p_i^2\right)^2 \\ &\leq p_1^4 \cdot \left(\sum_{i=2}^n p_i^2\right)^2 + p_1^2 \cdot \left(\sum_{i=2}^n p_i^2\right)^3 + \left(\sum_{i=2}^n p_i^2\right)^4 \\ &= \mathcal{O}\left(p_1^4 \cdot \left(\sum_{i=2}^n p_i^2\right)^2\right), \end{aligned}$$

where we used the facts that $\sum_{i=1}^n p_i^4 \leq (\sum_{i=1}^n p_i^2)^2$ and the prerequisite $\sum_{i=1}^n p_i^2 = \mathcal{O}(p_1^2)$. If we plug this into Equation (4.8), we get

$$\sum_A \sum_{B: |F_A \cap F_B|=l} \Pr(X_A \wedge X_B) = \mathcal{O} \left(m^{8-l} \cdot C^4 \cdot p_1^4 \cdot \left(\sum_{i=2}^n p_i^2 \right)^2 \cdot \sum_{\substack{B: \\ |F_A \cap F_B|=l}} \prod_{c \in F_B \setminus F_A} \Pr(c) \right). \quad (4.9)$$

Now we consider the cases $l \in \{1, 2, 3\}$. For $l = 1$ we know one clause which contains the central and one of the non-central variables. Thus, it holds that

$$\begin{aligned} \sum_{B: |F_A \cap F_B|=1} \prod_{c \in F_B \setminus F_A} \Pr(c) &\leq \left(\frac{C}{2} \right)^3 \cdot \sum_{x \in (S \cup \{i\})} p_x^3 \cdot \left(\sum_{y \in (S \cup \{i\}) \setminus \{x\}} p_y \cdot \left(\sum_{z \in [n] \setminus \{x, y\}} p_z^2 \right) \right) \\ &= \mathcal{O} \left(C^3 \cdot p_1^3 \cdot p_2 \left(\sum_{i=2}^n p_i^2 \right) \right), \end{aligned}$$

where the last line can be derived again by considering the possible cases for p_1 . Together with Equation (4.9), it now holds that

$$\begin{aligned} \sum_A \sum_{B: |F_A \cap F_B|=1} \Pr(X_A \wedge X_B) \\ = \mathcal{O} \left(m^7 \cdot C^7 \cdot p_1^7 \cdot p_2 \cdot \left(\sum_{i=2}^n p_i^2 \right)^3 \right) = o(1) = o(\mathbb{E}[X_2]^2), \quad (4.10) \end{aligned}$$

due to the choice of m and the prerequisite $p_2^2 = o(\sum_{i=2}^n p_i^2)$.

For $l = 2$ there can be two cases happening. Either all three variables appear in the two clauses or only two do. In the first case, one variable from $S \cup \{i\}$ appears in B twice as the center, while the other two appear only once. In the second case, one variable from $S \cup \{i\}$ appears in B twice again as the center and one new variable appears twice.

$$\sum_{B: |F_A \cap F_B|=2} \prod_{c \in F_B \setminus F_A} \Pr(c) = \mathcal{O} \left(C^2 \cdot p_1^2 \cdot p_2^2 + C^2 \cdot p_1^2 \cdot \left(\sum_{i=2}^n p_i^2 \right) \right),$$

Again with Equation (4.9), it holds that

$$\begin{aligned} \sum_A \sum_{B: |F_A \cap F_B|=2} \Pr(X_A \wedge X_B) \\ \leq \mathcal{O} \left(m^6 \cdot C^6 \cdot p_1^6 \cdot \left(\sum_{i=2}^n p_i^2 \right)^3 \right) = \mathcal{O}(1) = \mathcal{O}(\mathbb{E}[X_2]^2), \quad (4.11) \end{aligned}$$

where we used our choice of m again.

The last case is $l = 3$. This case can not happen, since the 3 clauses for B already fully determine the last clause, which also has to align with one of A , i. e. we do not have any degree of freedom to make $F_A \neq F_B$.

Putting equations (4.7), (4.10), and (4.11) together, establishes

$$\sum_A \sum_{B: B \sim A} \Pr(X_A \wedge X_B) = \mathcal{O}(\mathbb{E}[X_2]^2).$$

Together with Equation (4.6), this gives us

$$\mathbb{E}[X_2^2] = \sum_A \left(\sum_{B: B \not\sim A} \Pr(X_A \wedge X_B) + \sum_{B: B \sim A} \Pr(X_A \wedge X_B) \right) = \mathcal{O}(\mathbb{E}[X_2]^2)$$

and implies

$$\Pr(X_2 > 0) \geq \frac{\mathbb{E}[X_2]^2}{\mathbb{E}[X_2^2]} = \Omega(1).$$

as desired. \square

The following lemma complements the former one, showing that above that regime of m random formulas are unsatisfiable with probability $1 - o(1)$.

Lemma 4.5. *Consider a non-uniform random 2-SAT formula Φ with $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = o(\sum_{i=2}^n p_i^2)$. Then Φ is unsatisfiable with probability $1 - o(1)$ for $m = \omega\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right)$.*

Proof. We will show the result for $m = \omega\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right) \cap o(q_{\max}^{-1})$. For any $m = \Omega(q_{\max}^{-1})$ it follows by the fact that the probability that Φ is unsatisfiable is non-decreasing in m . However, the proof follows the same lines as the one for Lemma 4.4: We use the second moment method, but this time we want to show $\mathbb{E}[X_2^2] = (1 + o(1)) \cdot \mathbb{E}[X_2]^2$ in order to achieve

$$\Pr(X_2 > 0) \geq \frac{\mathbb{E}[X_2]^2}{\mathbb{E}[X_2^2]} = 1 - o(1).$$

Again, we look at the different parts of the following equation's right-hand side

$$\mathbb{E}[X_2^2] = \sum_A \sum_B \Pr(X_A \wedge X_B) = \sum_A \left(\sum_{B: B \not\sim A} \Pr(X_A \wedge X_B) + \sum_{B: B \sim A} \Pr(X_A \wedge X_B) \right).$$

Since our prerequisites ensure $m \cdot q_{\max} = o(1)$,

$$\sum_A \sum_{B: B \not\sim A} \Pr(X_A \wedge X_B) \leq (1 + o(1)) \mathbb{E}[X_2]^2. \quad (4.12)$$

still holds.

Again, we turn to the case $B \sim A$ and let $l = |F_A \cap F_B|$. Now we want to show that

$$\sum_A \sum_{B: B \sim A} \Pr(X_A \wedge X_B) = \sum_A \sum_{B: |F_A \cap F_B|=l} \Pr(X_A \wedge X_B) = o\left(\mathbb{E}[X_2]^2\right).$$

For $l = 4$ it holds that

$$\sum_A \sum_{B: |F_A \cap F_B|=4} \Pr(X_A \wedge X_B) = 4 \cdot \mathbb{E}[X_2] = o\left(\mathbb{E}[X_2]^2\right), \quad (4.13)$$

since now $\mathbb{E}[X_2] = \omega(1)$ due to Lemma 4.3. For $l = 1$ it still holds that

$$\sum_A \sum_{B: |F_A \cap F_B|=1} \Pr(X_A \wedge X_B) = \mathcal{O}\left(m^7 \cdot C^7 \cdot p_1^7 \cdot p_2 \cdot \left(\sum_{i=2}^n p_i^2\right)^3\right).$$

From Lemma 4.3 we know that

$$\mathbb{E}[X_2] \geq (1 - o(1)) \cdot m^4 \cdot C^4 \cdot p_1^4 \cdot \left(\sum_{i=1}^n p_i^2\right)^2 = \omega(1)$$

in our context. With $p_2^2 = o(\sum_{i=2}^n p_i^2)$ this implies

$$\sum_A \sum_{B: |F_A \cap F_B|=1} \Pr(X_A \wedge X_B) = o\left(m^7 \cdot C^7 \cdot p_1^7 \cdot \left(\sum_{i=2}^n p_i^2\right)^{7/2}\right) = o\left(\mathbb{E}[X_2]^2\right)$$

as desired. For $l = 2$ we still get

$$\sum_A \sum_{B: |F_A \cap F_B|=2} \Pr(X_A \wedge X_B) = \mathcal{O}\left(m^6 \cdot C^6 \cdot p_1^6 \cdot \left(\sum_{i=2}^n p_i^2\right)^3\right) = o\left(\mathbb{E}[X_2]^2\right).$$

Since the case $l = 3$ cannot happen, this already establishes $\sum_A \sum_{B: |F_A \cap F_B|=l} \Pr(X_A \wedge X_B) = o\left(\mathbb{E}[X_2]^2\right)$ as desired. \square

The former two lemmas together with Lemma 3.2 establish that in the case of $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$ and $p_2^2 = o(\sum_{i=2}^n p_i^2)$ the asymptotic threshold is at $m = \Theta\left(\left(C \cdot p_1 (\sum_{i=2}^n p_i^2)^{1/2}\right)^{-1}\right)$ and that it is coarse.

We now turn to the case $p_1^2 = \Theta(\sum_{i=1}^n p_i^2)$. Again, we have to consider different possibilities for the shared clauses of snakes A and B to influence $\Pr(X_A \wedge X_B)$. In the proofs of the former case this was rather easy, since we only considered the smallest possible snakes of size 3. Now the distinction becomes a bit more difficult. We will distinguish several cases: If the number of shared clauses is at least $t - 1$ then $\Pr(X_A \wedge X_B)$ is by roughly a factor of $(1 + \varepsilon)^t$ smaller than $\mathbb{E}[X_t]^2$. If the shared clauses form at least two connected sub-formulas, then there are enough variable appearances pre-defined for B to make $\Pr(X_A \wedge X_B)$ sufficiently small. The last case is that there is only one connected sub-formula, which is a lot smaller than $t - 1$. In that case we have to carefully consider what happens to the central variable from B , since this variable appears most times in B and the many appearances take degrees of freedom away from other variables, therefore making $\Pr(X_A \wedge X_B)$ small.

Lemma 4.6. *Consider a non-uniform random 2-SAT formula Φ with $p_1^2 = o(\sum_{i=1}^n p_i^2)$. Then Φ is unsatisfiable with probability $1 - o(1)$ for $m > (1 + \varepsilon) \cdot (\sum_{i=1}^n p_i^2)^{-1}$, where $\varepsilon > 0$ is a constant.*

Proof. Again, we utilize the second moment method. We want to show that F_A for a snake A of size t appears in Φ with probability $1 - o(1)$, i. e. Φ is almost surely unsatisfiable. This will hold for some $t = o\left(\sqrt{f(n)}\right) \cap \omega(\log f(n))$, where $f(n) = (\sum_{i=1}^n p_i^2) / (p_1^2)$. Thus, we choose $t = \Theta(\log^2 f(n))$. We will later see why we chose t this way. Again, we define X_A as an indicator variable for the event that the formula F_A associated with snake A appears exactly once in Φ and

$$X_t = \sum_{\text{snake } A \text{ of size } t} X_A.$$

As in the proof of Lemma 4.5 we want to show $\mathbb{E}[X_t^2] \leq (1 + o(1)) \cdot \mathbb{E}[X_t]^2$, giving us the desired result. We again split the expected value into two sums

$$\mathbb{E}[X_t^2] = \sum_A \sum_B \Pr(X_A \wedge X_B) = \sum_{B: B \approx A} \Pr(X_A \wedge X_B) + \sum_{B: B \sim A} \Pr(X_A \wedge X_B),$$

where $B \sim A$ denotes $F_A \cap F_B \neq \emptyset$. We will now consider the parts over $B \approx A$ and $B \sim A$ separately, starting with $B \approx A$.

As in the proof of Lemma 4.5, we want to show that

$$\sum_A \sum_{B: B \approx A} \Pr(X_A \wedge X_B) = (1 + o(1)) \cdot \mathbb{E}[X_t]^2. \quad (4.14)$$

It holds that

$$\Pr(X_A \wedge X_B) = \binom{m}{4t} \cdot (4t)! \cdot \left(\prod_{c \in F_A} \Pr(c)\right) \cdot \left(\prod_{c \in F_B} \Pr(c)\right) \cdot \left(1 - \sum_{c \in F_A \cup F_B} \Pr(c)\right)^{m-4t},$$

while

$$\Pr(X_A) = \binom{m}{2t} \cdot (2t)! \cdot \left(\prod_{c \in F_A} \Pr(c)\right) \cdot \left(1 - \sum_{c \in F_A} \Pr(c)\right)^{m-2t}.$$

This already gives us

$$\Pr(X_A \wedge X_B) \leq \Pr(X_A) \cdot \Pr(X_B) \frac{\left(1 - \sum_{c \in F_A \cup F_B} \Pr(c)\right)^{m-4t}}{\left(1 - \sum_{c \in F_A} \Pr(c)\right)^{m-2t} \left(1 - \sum_{c \in F_B} \Pr(c)\right)^{m-2t}},$$

since $\binom{m}{4t} \cdot (4t)! \leq \left(\binom{m}{2t} \cdot (2t)!\right)^2$. W.l.o.g. $\sum_{c \in F_B} \Pr(c) \leq \sum_{c \in F_A} \Pr(c)$. Now it holds that

$$\begin{aligned}
& \left(1 - \sum_{c \in F_A \cup F_B} \Pr(c)\right)^{m-4t} \left(\frac{1}{1 - \sum_{c \in F_A} \Pr(c)}\right)^{m-2t} \left(\frac{1}{1 - \sum_{c \in F_B} \Pr(c)}\right)^{m-2t} \\
&= \left(1 - \frac{\sum_{c \in F_B} \Pr(c)}{1 - \sum_{c \in F_A} \Pr(c)}\right)^{m-4t} \left(\frac{1}{1 - \sum_{c \in F_A} \Pr(c)}\right)^{2t} \left(\frac{1}{1 - \sum_{c \in F_B} \Pr(c)}\right)^{m-2t} \\
&\leq \exp\left(- (m-4t) \cdot \frac{\sum_{c \in F_B} \Pr(c)}{1 - \sum_{c \in F_A} \Pr(c)} + 2t \cdot \frac{\sum_{c \in F_A} \Pr(c)}{1 - \sum_{c \in F_A} \Pr(c)} + (m-2t) \cdot \frac{\sum_{c \in F_B} \Pr(c)}{1 - \sum_{c \in F_B} \Pr(c)}\right) \\
&\leq \exp\left(2t \cdot \frac{\sum_{c \in F_A} \Pr(c)}{1 - \sum_{c \in F_A} \Pr(c)} + 2t \cdot \frac{\sum_{c \in F_B} \Pr(c)}{1 - \sum_{c \in F_B} \Pr(c)}\right) \\
&\leq \exp\left(4t \cdot \frac{2t \cdot q_{\max}}{1 - 2t \cdot q_{\max}}\right),
\end{aligned}$$

where the second-to-last line followed with $\sum_{c \in F_A} \Pr(c) \leq \sum_{c \in F_A} \Pr(c)$ and q_{\max} is the maximum clause probability. This expression is $1 + o(1)$, since $t^2 \cdot q_{\max} = o(C \cdot f(n) \cdot p_1^2) = o(\sum_{i=1}^n p_i^2) = o(1)$. We now get

$$\Pr(X_A \wedge X_B) = (1 + o(1)) \cdot \Pr(X_A) \cdot \Pr(X_B)$$

for $A \approx B$ and thus

$$\sum_A \sum_{B: B \approx A} \Pr(X_A \wedge X_B) \leq (1 + o(1)) \cdot \sum_A \sum_{B: B \approx A} \Pr(X_A) \Pr(X_B) \leq (1 + o(1)) \cdot \mathbb{E}[X_2]^2.$$

Second, we look at snakes $B \sim A$. For those we want to show

$$\sum_{B: B \sim A} \Pr(X_A \wedge X_B) = o\left(\mathbb{E}[X_t]^2\right). \quad (4.15)$$

This now becomes a bit more complicated than in the case of $t = 2$, since we can not always surely say how many variables are predefined by shared clauses in snake B . As before, we are now classifying snakes $B \sim A$ according to the number $l = |F_A \cap F_B|$ of shared clauses, but also according to the number j of nodes in the variable-variable incidence graph $G_{F_A \cap F_B}$. Note that actually, the number of variables that F_A and F_B have in common (regardless of signs) could be greater! In fact, they could share all their variables without having a single clause in common. However, right now we are only interested in ways to incorporate clauses from F_A as common clauses into F_B . To that end, we only need to consider the variables from these clauses as shared variables. Suppose now that snake A is fixed. We now know that there are $2t - 1 - j$ free variables in B , i. e. variables which are not predetermined by shared clauses. Furthermore we can give an upper bound on the number c of connected components of $G_{F_A \cap F_B}$. It is easy to see that $c \leq j - l$ for $l < t$ ($G_{F_A \cap F_B}$ is a forest), $c \leq j - l + 1$ for $t \leq l < 2t$ (we could create a cycle), and $c = j - l + 2$ for $l = 2t$ ($F_A = F_B$). Fixing l and j it holds that

$$\begin{aligned}
& \sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})|=j}} \Pr(X_A \wedge X_B) \\
&\leq \binom{m}{4t-l} \cdot (4t-l)! \cdot \left(\frac{C}{2}\right)^{4t-l} \cdot 2^{2t-1} \cdot (2t-2)! \left(\sum_{\substack{S_A \subseteq [n]: \\ |S_A|=2t-2}} \prod_{x \in S_A} p(x)^2 \cdot \sum_{y \in [n]} p(y)^4 \right) \\
&\cdot 4 \left(\binom{2t+2}{2(j-l)+2} \right)^2 \cdot c! \cdot 2^c \cdot 2t \cdot (2t-1-j)! \cdot 2^{2t-1-j} \cdot \left(\sum_{\substack{S_B \subseteq [n]: \\ |S_B|=2t-1-j}} \prod_{x \in S_B} p(x)^2 \right) \\
&\cdot p_1^{2(j-l+1)} \cdot \left(1 - \sum_{c \in F_A \cup F_B} P(c)\right)^{m-(4t-l)}. \quad (4.16)
\end{aligned}$$

Before we upper bound this expression even further, let us explain where it comes from. There are $\binom{m}{4t-l} \cdot (4t-l)!$ positions for the $4t-l$ clauses of $F_A \cup F_B$ in the m -clause formula Φ . There are at

most $2^{2t-2} \cdot (2t-2)!$ possibilities of forming different snakes (signs and positions) from the $2t-2$ variables excluding $|w_t|$ and two possible signs for $|w_t|$. In snake A each variable appears exactly twice, except for $|w_t|$, which appears four times. Now we want to count the ways of mapping $G_{F_A \cap F_B}$ to G_{F_A} and G_{F_B} respectively. Following the argumentation from [11] we can see that there are $2^{\binom{2t+2}{2j-2l+2}}$ possible mappings for G_{F_A} and G_{F_B} respectively. These mappings fix the shared clauses we choose from A as well as the positions where shared clauses can appear in B , but not where exactly which clause will appear. We know that $G_{F_A \cap F_B}$ contains c connected components. If they are of same length, they can be interchanged in $c!$ ways. Furthermore, each component might be flipped, i. e. the sign of every literal in the component and their order in B can be inverted. For components which are paths, this does not change the set of shared clauses they originate from. Nevertheless, there is still the possibility of having one component which is not a path. For this component there are at most $2t$ ways of mapping it onto its counterpart (if it is a cycle) due to [11]. Now we know the shared clauses from F_A and the exact position of these clauses in F_B as well as positions reserved for non-determined variables in snake B . The remaining $2t-1-j$ non-determined variables from B can be chosen arbitrarily. Also, there are $2^{2t-1-j} \cdot (2t-1-j)!$ possibilities for them to fill out the blanks of snake B . The remaining at most $2(j-l+1)$ appearances of variables in F_B are determined by the previous choices and give an additional factor of at most $p_1^{2(j-l+1)}$. Note that the case that one of our free variables in B is a central variable is also captured by this upper bound, since $\sum_{i=1}^n p_i^4 \leq p_1^2 \cdot \sum_{i=1}^n p_i^2$. The other $m-(4t-l)$ clauses of F are supposed to be different from those in $F_A \cup F_B$, so that both F_A and F_B appear exactly once.

Now we want to simplify that expression a bit. It holds that $(1 - \sum_{C \in F_A \cup F_B} P(C))^{m-(4t-l)} \leq 1$ and that

$$C^{4t-l} \leq \left(1 + \frac{\sum_{i=1}^n p_i^2}{1 - \sum_{i=1}^n p_i^2}\right)^{4t} \leq \exp\left(4t \cdot \frac{\sum_{i=1}^n p_i^2}{1 - \sum_{i=1}^n p_i^2}\right) = 1 + o(1),$$

since $t \cdot \sum_{i=1}^n p_i^2 = o(f(n)^{1/2} \cdot \sum_{i=1}^n p_i^2) = o(f(n)^{3/2} \cdot p_1^2) = o(f(n)^{-1/2}) = o(1)$ due to $p_1 \leq f(n)^{-1}$. Again,

$$\sum_{\substack{S \subseteq [n]: \\ |S|=x}} \prod_{s \in S} p(s)^2 \leq \frac{1}{x!} \left(\sum_{i=1}^n p_i^2\right)^x.$$

This step also cancels out the factors $(2t-2)!$ and $(2t-1-j)!$. Also, all factors of 2 that appear cancel out with $c \leq j-l+2$. We will also use the following estimation

$$\left(\binom{2t+2}{2(j-l)+2}\right)^2 \cdot c! \leq \frac{(2t+2)^{4(j-l+1)}}{(2(j-l+1)!)^2} \cdot (j-l+2)! \leq (2t+2)^{4(j-l+1)}.$$

Plugging everything back into Equation (4.16) we get

$$\begin{aligned} & \sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})|=j}} \Pr(X_A \wedge X_B) \\ & \leq 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^{5(j-l+1)} \cdot \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=1}^n p_i^2\right)^{4t-j-3} \cdot p_1^{2(j-l+1)} \end{aligned} \quad (4.17)$$

Remember that due to Lemma 4.2

$$\mathbb{E}[X_t]^2 \geq (1 - o(1)) \cdot m^{4t} \left(\sum_{i=1}^n p_i^4\right)^2 \cdot \left(\sum_{i=1}^n p_i^2\right)^{4t-4}.$$

We will distinguish three cases now, depending on the value of $j-l$. First $j-l=0$, then $j-l \geq 2$ and finally $j-l=1$. For each of these cases we want to show

$$\sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})|=j}} \Pr(X_A \wedge X_B) = o\left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)}\right).$$

Since $1 \leq l \leq 2t$ and $2 \leq j \leq 2t - 1$, we will get an additional factor of $4t^2$ when summing over all snakes $A \sim B$. With our choice $t = \Theta(\log^2 f(n))$, this adds up to

$$\sum_{B: B \sim A} \Pr(X_A \wedge X_B) = o\left(\mathbb{E}[X_t]^2\right)$$

as desired.

Now let us consider the first case, $j = l$. This can only happen if $G_{F_A \cap F_B}$ contains a cycle, i. e. $l \geq t$. It now holds that

$$\begin{aligned} & \sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})|=l}} \Pr(X_A \wedge X_B) \\ & \leq 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^5 \cdot \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=1}^n p_i^2\right)^{4t-l-3} \cdot p_1^2 \\ & = \mathcal{O}\left(t^5 \left(m \cdot \sum_{i=1}^n p_i^2\right)^{-l} \cdot \frac{p_1^2 \cdot \sum_{i=1}^n p_i^2}{\sum_{i=1}^n p_i^4} \cdot \mathbb{E}[X_t]^2\right) \\ & = \mathcal{O}\left(t^5 (1+\varepsilon)^{-t} \cdot f(n) \cdot \mathbb{E}[X_t]^2\right) = o\left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)}\right), \end{aligned}$$

due to $\sum_{i=1}^n p_i^4 \geq p_1^4$ and due to our choice $t = \Theta(\log^2 f(n))$.

The second case we consider is $j - l \geq 2$. It holds that

$$\begin{aligned} & \sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})| \geq l+2}} \Pr(X_A \wedge X_B) \\ & \leq 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^{5(j-l+1)} \cdot \left(\sum_{i=1}^n p_i^4\right) \cdot \left(\sum_{i=1}^n p_i^2\right)^{4t-j-3} \cdot p_1^{2(j-l+1)} \\ & = \mathcal{O}\left(t^{5(j-l+1)} \cdot m^{-l} \left(\sum_{i=1}^n p_i^2\right)^{-j+1} \cdot \left(\sum_{i=1}^n p_i^4\right)^{-1} \cdot p_1^{2(j-l+1)} \cdot \mathbb{E}[X_t]^2\right) \\ & = \mathcal{O}\left(t^{5(j-l+1)} \left(m \cdot \sum_{i=1}^n p_i^2\right)^{-l} \cdot \frac{p_1^{2(j-l+1)}}{p_1^4 \left(\sum_{i=1}^n p_i^2\right)^{j-l-1}} \cdot \mathbb{E}[X_t]^2\right) \\ & = \mathcal{O}\left(t^{5(j-l+1)} \cdot \frac{p_1^{2(j-l-1)}}{\left(\sum_{i=1}^n p_i^2\right)^{j-l-1}} \cdot \mathbb{E}[X_t]^2\right) \\ & = \mathcal{O}\left(t^{5(j-l+1)} f(n)^{-(j-l-1)} \cdot \mathbb{E}[X_t]^2\right) \\ & = \mathcal{O}\left(t^{10} \cdot \left(\frac{t^5}{f(n)}\right)^{j-l-1} \cdot \mathbb{E}[X_t]^2\right) = o\left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)}\right), \end{aligned}$$

since $j - l - 1 \geq 1$ and $t = \Theta(\log^2 f(n))$.

The last case we consider is $j - l = 1$. This happens if we either only have one connected component in $G_{F_A \cap F_B}$ that does not form a cycle or if $G_{F_A \cap F_B}$ contains a cycle and one other

connected component. In the latter case, we get

$$\begin{aligned}
& \sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})|=l+1 \\ \text{cycle in } G_{F_A \cap F_B}}} \Pr(X_A \wedge X_B) \\
& \leq 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^{10} \cdot \left(\sum_{i=1}^n p_i^4 \right) \cdot \left(\sum_{i=1}^n p_i^2 \right)^{4t-l-4} \cdot p_1^4 \\
& = \mathcal{O} \left(t^{10} \left(m \cdot \sum_{i=1}^n p_i^2 \right)^{-l} \cdot \frac{p_1^4}{\sum_{i=1}^n p_i^4} \cdot \mathbb{E}[X_t]^2 \right) \\
& = \mathcal{O} \left(t^{10} (1 + \varepsilon)^{-t} \cdot \mathbb{E}[X_t]^2 \right) = o \left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)} \right),
\end{aligned}$$

since a cycle can only exist for $l \geq t$ and since we choose $t = \Theta(\log^2 f(n))$.

If $G_{F_A \cap F_B}$ that does not form a cycle, we have to look a bit more closely now, since we cannot guarantee a large enough t to make the expression sufficiently small. Instead, we will consider different cases for the central variable in B . First, we assume that the central variable is a free variable. Then, we actually get

$$\begin{aligned}
& 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^{10} \cdot \left(\sum_{i=1}^n p_i^4 \right)^2 \cdot \left(\sum_{i=1}^n p_i^2 \right)^{4t-l-5} \cdot p_1^2 \\
& = \mathcal{O} \left(t^{10} \left(m \cdot \sum_{i=1}^n p_i^2 \right)^{-l} \cdot \frac{p_1^2}{\sum_{i=1}^n p_i^2} \cdot \mathbb{E}[X_t]^2 \right) \\
& = \mathcal{O} \left(t^{10} f(n)^{-1} \cdot \mathbb{E}[X_t]^2 \right) = o \left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)} \right),
\end{aligned}$$

since now we have a second variable that we can choose freely and which appears at least 4 times.

Now we assume that the central variable in B is not free. What could happen? It could coincide with a non-central variable from A or with the central variable from A . Also, the central variable could already appear once or twice in shared clauses in the first and one to four times in the second case.

Let us start with the case that it coincides with a non-central variable in A . Then, one of the variables that appears twice in A appears an additional (not in shared clauses) 2 or 3 times as the central node in B , depending on the number of shared clauses it already appears in. So, in total it either appears 4 times or 5 times, replacing one appearance of a variable that appears twice in A and 2 resp. 3 appearances of unfree variables in B . Since $\sum_{i=1}^n p_i^5 \leq p_1 \sum_{i=1}^n p_i^4$, the former case gives us an upper bound. We get at most

$$\begin{aligned}
& 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^{10} \cdot \left(\sum_{i=1}^n p_i^4 \right)^2 \cdot \left(\sum_{i=1}^n p_i^2 \right)^{4t-l-5} \cdot p_1^2 \\
& = \mathcal{O} \left(t^{10} \left(m \cdot \sum_{i=1}^n p_i^2 \right)^{-l} \cdot \frac{p_1^2}{\sum_{i=1}^n p_i^2} \cdot \mathbb{E}[X_t]^2 \right) \\
& = \mathcal{O} \left(t^{10} f(n)^{-1} \cdot \mathbb{E}[X_t]^2 \right) = o \left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)} \right).
\end{aligned}$$

The other case is that it coincides with the central variable from A . Then, the variable that appears 4 times in A might appear 0 to 3 additional times (not in shared clauses) in B , depending on the number of shared clauses it already appears in. It cannot appear an additional 4 times, since the central variable of A must appear in a shared clause at least once for the variable to not be free. However, this means that all our unfree variables actually belong to distinct variables that appear twice in A and an additional time in B . If some of them belonged to the same variable, this would again imply $l \geq t-1$ and we could handle this case by having

a large enough l . Let $x \in \{0, 1, 2, 3\}$ be the number of times that the central variable from A appears additionally in B . We now get

$$\begin{aligned} & 4 \cdot (1 + o(1)) \cdot m^{4t-l} \cdot (2t+2)^{10} \cdot \left(\sum_{i=1}^n p_i^{4+x} \right) \cdot \left(\sum_{i=1}^n p_i^2 \right)^{4t-l-4-(4-x)} \cdot \left(\sum_{i=1}^n p_i^3 \right)^{4-x} \\ &= \mathcal{O} \left(t^{10} \left(m \cdot \sum_{i=1}^n p_i^2 \right)^{-l} \cdot \frac{(\sum_{i=1}^n p_i^{4+x}) \cdot (\sum_{i=1}^n p_i^3)^{4-x}}{(\sum_{i=1}^n p_i^4)^2 (\sum_{i=1}^n p_i^2)^{4-x}} \cdot \mathbb{E}[X_t]^2 \right) \\ &= \mathcal{O} \left(t^{10} \cdot \frac{(\sum_{i=1}^n p_i^{4+x}) \cdot (\sum_{i=1}^n p_i^3)^{4-x}}{(\sum_{i=1}^n p_i^4)^2 (\sum_{i=1}^n p_i^2)^{4-x}} \cdot \mathbb{E}[X_t]^2 \right). \end{aligned}$$

It remains to show that

$$t^{10} \cdot \frac{(\sum_{i=1}^n p_i^{4+x}) \cdot (\sum_{i=1}^n p_i^3)^{4-x}}{(\sum_{i=1}^n p_i^4)^2 (\sum_{i=1}^n p_i^2)^{4-x}} = o \left(\frac{1}{\log^4 f(n)} \right).$$

In order to do so, consider p_1, p_2, \dots, p_n . We now split the probabilities into those with $p_i \geq p_1 / \log^y f(n)$ and those with $p_i < p_1 / \log^y f(n)$, where $y \in \mathbb{N}$ will be determined later. Now let $N_{\max} = \{i \in [n] \mid p_i \geq p_1 / \log^y f(n)\}$. It holds that

$$N_{\max} \cdot \left(\frac{p_1}{\log^y f(n)} \right)^2 + \sum_{i \in [n]: p_i < p_1 / \log^y f(n)} p_i^2 \leq \sum_{i=1}^n p_i^2 = f(n) \cdot p_1^2.$$

This implies $N_{\max} = \mathcal{O}(f(n) \cdot \log^{2y} f(n))$. We now distinguish two cases: $N_{\max} \geq f(n)^{2/3} \cdot \log^{4y/3} f(n)$ and $N_{\max} < f(n)^{2/3} \cdot \log^{4y/3} f(n)$.

Now assume the first case, $N_{\max} \geq f(n)^{2/3} \cdot \log^{4y/3} f(n)$. It holds that

$$\sum_{i=1}^n p_i^4 \geq f(n)^{2/3} \cdot \log^{4y/3} f(n) \cdot \left(\frac{p_1}{\log^y f(n)} \right)^4 = \frac{f(n)^{2/3}}{\log^{(8/3)y} f(n)} \cdot p_1^4.$$

This implies

$$\begin{aligned} t^{10} \cdot \frac{(\sum_{i=1}^n p_i^{4+x}) \cdot (\sum_{i=1}^n p_i^3)^{4-x}}{(\sum_{i=1}^n p_i^4)^2 (\sum_{i=1}^n p_i^2)^{4-x}} &\leq t^{10} \cdot \frac{f(n) \cdot p_1^{4+x} \cdot (f(n) \cdot p_1^3)^{4-x} \cdot \log^{(16/3)y} f(n)}{f(n)^{4/3} \cdot p_1^8 \cdot (f(n) \cdot p_1^2)^{4-x}} \\ &= t^{10} \cdot f(n)^{-1/3} \cdot \log^{(16/3)y} f(n) = o \left(\frac{1}{\log^4 f(n)} \right) \end{aligned}$$

as desired, since $\sum_{i=1}^n p_i^x \leq p_1^{x-2} \cdot \sum_{i=1}^n p_i^2 = f(n) \cdot p_1^x$ for $x \in \mathbb{N}$ with $x \geq 3$ and $t = \Theta(\log^2 f(n))$.

Now assume $N_{\max} < f(n)^{2/3} \cdot \log^{4y/3} f(n)$. It holds that

$$\begin{aligned} \sum_{i=1}^n p_i^3 &< f(n)^{2/3} \cdot \log^{4y/3} f(n) \cdot p_1^3 + \frac{p_1}{\log^y f(n)} \cdot \sum_{i \in [n]: p_i < p_1 / \log^y f(n)} p_i^2 \\ &\leq f(n)^{2/3} \cdot \log^{4y/3} f(n) \cdot p_1^3 + \frac{f(n)}{\log^y f(n)} \cdot p_1^3 \\ &= \mathcal{O} \left(\frac{f(n)}{\log^y f(n)} \cdot p_1^3 \right). \end{aligned}$$

This readily implies

$$\begin{aligned} t^{10} \cdot \frac{(\sum_{i=1}^n p_i^{4+x}) \cdot (\sum_{i=1}^n p_i^3)^{4-x}}{(\sum_{i=1}^n p_i^4)^2 (\sum_{i=1}^n p_i^2)^{4-x}} &= \mathcal{O} \left(t^{10} \cdot \frac{(\sum_{i=1}^n p_i^4) \cdot p_1^x \cdot \left(\frac{f(n)}{\log^y f(n)} \cdot p_1^3 \right)^{4-x}}{p_1^4 \cdot (\sum_{i=1}^n p_i^4) \cdot (f(n) \cdot p_1^2)^{4-x}} \right) \\ &= \mathcal{O} \left(\frac{t^{10}}{\log^{y(4-x)} f(n)} \right) = o \left(\frac{1}{\log^4 f(n)} \right) \end{aligned}$$

where we used $\sum_{i=1}^n p_i^{4+x} \leq p_1^x \cdot \sum_{i=1}^n p_i^4$ and $\sum_{i=1}^n p_i^4 \geq p_1^4$ in the first line and $t = \Theta(\log^2 f(n))$ with $x \leq 3$ and $y = 25$ in the last line.

Finally, we took care of all the cases for $j - l = 1$ and showed

$$\sum_{\substack{\text{snakes } A, B: \\ |E(G_{F_A \cap F_B})|=l, |V(G_{F_A \cap F_B})|=l+1}} \Pr(X_A \wedge X_B) = o\left(\frac{\mathbb{E}[X_t]^2}{\log^4 f(n)}\right)$$

as desired. This implies

$$\sum_{B: B \sim A} \Pr(X_A \wedge X_B) = o\left(\mathbb{E}[X_t]^2\right)$$

and concludes the proof. \square

Lemma 4.6 and Lemma 3.1 now establish the existence of a sharp threshold at $m = \left(\sum_{i=1}^n p_i^2\right)^{-1}$.

Now we still have to consider the case $p_1^2 = \Theta\left(\sum_{i=1}^n p_i^2\right)$ and $p_2^2 = \Theta\left(\sum_{i=2}^n p_i^2\right)$. In the following lemma, we give a lower bound on the probability to be unsatisfiable by showing the existence of an unsatisfiable sub-formula consisting only of the two most-probable variables. The lemma generally holds for $k \geq 2$, but it especially serves our purpose of considering the remaining case.

Lemma 4.7. *Consider a non-uniform random k -SAT formula Φ with $q_{\max} = o(1)$. Then Φ is unsatisfiable with probability at least*

$$(1 - \exp(-q_{\max} \cdot m))^{2^k} - q_{\max}^2 \cdot 2^{2k} \cdot m \cdot (1 + \exp(-q_{\max} \cdot m))^{2^k}.$$

Proof. Let c be the clause with maximum probability. Since the signs of literals are chosen with probability $1/2$ independently at random, it holds that each clause with the same variables as c has the same probability. Our lower bound is now just a lower bound on the probability of having each of the 2^k clauses with these variables, which constitute an unsatisfiable sub-formula. Let us enumerate the different clauses c_1, \dots, c_{2^k} with variables X_1, \dots, X_k in an arbitrary order. Now let $\bar{\mathcal{E}}_j$ denote the event that c_j is *not* appearing in Φ and let $\bar{\mathcal{E}} = \cup_{j=1}^{2^k} \bar{\mathcal{E}}_j$ denote the event that at least one of these clauses does not appear. Due to the principle of inclusion and exclusion it holds that

$$\Pr(\bar{\mathcal{E}}) = \sum_{l=1}^{2^k} (-1)^{l+1} \sum_{J \subseteq [2^k]: |J|=l} \Pr\left(\bigcap_{j \in J} \bar{\mathcal{E}}_j\right) = \sum_{l=1}^{2^k} (-1)^{l+1} \binom{2^k}{l} \cdot (1 - l \cdot q_{\max})^m,$$

because the clauses c_1, \dots, c_{2^k} have the same probability q_{\max} of appearing and all clauses are drawn independently at random.

It now holds that

$$\begin{aligned} \Pr(\Phi \text{ unsat}) &\geq 1 - \left(\sum_{l=1}^{2^k} \binom{2^k}{l} \cdot (-1)^l \cdot (1 - l \cdot q_{\max})^m\right) \\ &= \sum_{l=0}^{2^k} \binom{2^k}{l} \cdot (-1)^l \cdot (1 - l \cdot q_{\max})^m. \end{aligned}$$

We can now estimate

$$-(1 - q_{\max} \cdot l)^m \geq -\exp(-q_{\max} \cdot l \cdot m)$$

and, due to [37, Proposition B.3],

$$\begin{aligned} (1 - q_{\max} \cdot l)^m &\geq \exp(-q_{\max} \cdot l \cdot m) \cdot (1 - q_{\max}^2 \cdot l^2 \cdot m) \\ &\geq \exp(-q_{\max} \cdot l \cdot m) \cdot (1 - q_{\max}^2 \cdot 2^{2k} \cdot m). \end{aligned}$$

In total, we get

$$\begin{aligned}
& \Pr(\Phi \text{ unsat}) \\
& \geq \sum_{l=0}^{2^k} \left(\binom{2^k}{l} \cdot (-1)^l \cdot \exp(-q_{\max} \cdot l \cdot m) - \binom{2^k}{l} \cdot q_{\max}^2 \cdot 2^{2k} \cdot m \cdot \exp(-q_{\max} \cdot l \cdot m) \right) \\
& = (1 - \exp(-q_{\max} \cdot m))^{2^k} - q_{\max}^2 \cdot 2^{2k} \cdot m \cdot (1 + \exp(-q_{\max} \cdot m))^{2^k}.
\end{aligned}$$

□

The former lemma now yields the following corollary.

Corollary 4.1. *Consider a non-uniform random k -SAT formula Φ with $q_{\max} = o(1)$. Then*

1. $\Pr(\Phi \text{ unsatisfiable}) = \Omega(1)$ for $m = \Theta(q_{\max}^{-1})$ and
2. $\Pr(\Phi \text{ unsatisfiable}) = 1 - o(1)$ for $m = \omega(q_{\max}^{-1})$.

In the second case the result follows from Lemma 4.7 for $m = \omega(q_{\max}^{-1}) \cap o(q_{\max}^{-2})$ and by monotonicity of the satisfiability probability in m . This corollary together with Lemma 3.2 establishes the existence of a coarse threshold at $m = \Theta\left(\left(C \cdot p_1 \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)^{-1}\right) = \Theta(q_{\max}^{-1})$ for non-uniform random 2-SAT with $p_1^2 = \Theta\left(\sum_{i=1}^n p_i^2\right)$, $p_2^2 = \Theta\left(\sum_{i=2}^n p_i^2\right)$.

5 Constant Clause Probabilities

We assumed $q_{\max} = o(1)$ throughout the paper. For the sake of completeness we still have to take care of the case $q_{\max} = \Theta(1)$. It is easy to see that for $\Phi \sim \mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ and a constant $m \geq 4$ it holds that $\Pr(\Phi \text{ unsatisfiable}) \geq q_{\max}^m$, since this is the probability of an unsatisfiable instance, where the most probable clause appears with all four combinations of signs and then one of these clauses appears an additional $m - 4$ times. Similarly, $\Pr(\Phi \text{ satisfiable}) \geq q_{\max}^m$, as this is the probability of a satisfiable instance, where the same most probable clause appears m times with the same sign. Since $0 < q_{\max} \leq 1/4$ is a constant, the probability is a constant bounded away from zero and one. It remains to show that Φ is unsatisfiable with probability $1 - o(1)$ for $m = \omega(1)$. The following lemma establishes this. Again, this lemma also holds for $k \geq 2$ in general.

Lemma 5.1. *Consider a non-uniform random k -SAT formula Φ . Then Φ is unsatisfiable with probability at least*

$$2 - (1 + \exp(-q_{\max} \cdot m))^{2^k}.$$

Proof. As in Lemma 4.7, it holds that

$$\Pr(\Phi \text{ unsat}) \geq \sum_{l=0}^{2^k} \left(\binom{2^k}{l} (-1)^l (1 - l \cdot q_{\max})^m \right).$$

We can now estimate

$$\begin{aligned}
\sum_{l=0}^{2^k} \left(\binom{2^k}{l} (-1)^l (1 - l \cdot q_{\max})^m \right) & \geq 1 - \sum_{l=1}^{2^k} \left(\binom{2^k}{l} (1 - l \cdot q_{\max})^m \right) \\
& \geq 1 - \sum_{l=1}^{2^k} \left(\binom{2^k}{l} \exp\left(-m \cdot \frac{l \cdot q_{\max}}{1 - l \cdot q_{\max}}\right)^m \right) \\
& \geq 1 - \sum_{l=1}^{2^k} \left(\binom{2^k}{l} \exp(-m \cdot l \cdot q_{\max}) \right) \\
& = 2 - (1 + \exp(-m \cdot q_{\max}))^{2^k}
\end{aligned}$$

□

For $q_{\max} = \Theta(1)$ and $m = \omega(q_{\max}^{-1})$ this lemma implies $\Pr(\Phi \text{ unsatisfiable}) \geq 1 - o(1)$. All lemmas together now imply our main theorem.

Theorem 1.1. *Let $\mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ be the non-uniform random 2-SAT model with n variables, m clauses, and an ensemble of probability distributions $(\vec{p}_x)_{x \in \mathbb{N}}$. Let $\vec{p}_n = (p_1, p_2, \dots, p_n)$ be the n -th distribution from the ensemble. W.l.o.g. let $p_1 \geq p_2 \geq \dots \geq p_n$. If $p_1^2 = o(\sum_{i=1}^n p_i^2)$, then $\mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ has a sharp satisfiability threshold at $m = (\sum_{i=1}^n p_i^2)^{-1}$. Otherwise, $\mathcal{D}(n, 2, (\vec{p}_x)_{x \in \mathbb{N}}, m)$ has a coarse satisfiability threshold at $m = \Theta\left(\left(1 - \sum_{i=1}^n p_i^2\right) / \left(p_1 \cdot \left(\sum_{i=2}^n p_i^2\right)^{1/2}\right)\right)$.*

6 Example Applications of our Theorem

We will now show on some examples how our main theorem can be applied.

6.1 Uniform Distribution

The simplest distribution we can apply our theorem to is the uniform distribution, i.e. $\vec{p}_n = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$ for all $n \in \mathbb{N}$. It holds that $p_1^2 = \frac{1}{n^2}$ and $\sum_{i=1}^n p_i^2 = \frac{1}{n}$. Thus, Theorem 1.1 implies a sharp threshold at $m^*(n) = n$ for all $n \in \mathbb{N}$. This proves the satisfiability threshold conjecture for $k = 2$.

6.2 Power Law Distribution

Another ensemble of distributions we can choose are power-law distributions, i.e. we consider the power law random 2-SAT model introduced by Ansótegui et al. [3]. Thus, for a constant $\beta > 2$ we choose $\vec{p}_n = (p_1, p_2, \dots, p_n)$ with

$$p_i = \frac{(n/i)^{\frac{1}{\beta-1}}}{\left(\sum_{j=1}^n (n/j)^{\frac{1}{\beta-1}}\right)}.$$

It already holds that $p_1 \geq p_2 \geq \dots \geq p_n$. Now it is an easy exercise to show that

$$\left(\sum_{j=1}^n (n/j)^{\frac{1}{\beta-1}}\right) = (1 - o(1)) \cdot \frac{\beta - 1}{\beta - 2}.$$

Furthermore

$$p_1^2 = (1 \pm o(1)) \cdot \left(\frac{\beta - 1}{\beta - 2}\right)^2 \cdot n^{-2\frac{\beta-2}{\beta-1}}.$$

Finally, one can show that

$$\sum_{i=1}^n p_i^2 = \begin{cases} (1 \pm o(1)) \cdot \frac{(\beta-2)^2}{(\beta-3) \cdot (\beta-1)} \cdot n^{-2\frac{\beta-2}{\beta-1}} & \text{for } \beta < 3 \\ (1 \pm o(1)) \cdot \frac{1}{4} \cdot \frac{\ln n}{n} & \text{for } \beta = 3 \\ (1 \pm o(1)) \cdot \frac{(\beta-2)^2}{(\beta-3) \cdot (\beta-1)} \cdot n^{-1} & \text{for } \beta > 3. \end{cases}$$

Thus, applying our theorem we can see that for $\beta < 3$ there is a coarse threshold at $m = \Theta(n^{-2\frac{\beta-2}{\beta-1}})$, since $p_1^2 = \Theta(\sum_{i=1}^n p_i^2) = \Theta(n^{-2\frac{\beta-2}{\beta-1}})$ and $C = 1 + o(1)$. For $\beta = 3$ there is a sharp threshold at $4 \cdot \frac{n}{\ln n}$, since $p_1^2 = \Theta(n^{-1}) = o(\frac{\ln n}{n})$. Also, there is a sharp threshold at $\frac{(\beta-3) \cdot (\beta-1)}{(\beta-2)^2} \cdot n$ for $\beta > 3$, since $p_1^2 = \Theta\left(n^{-2\frac{\beta-2}{\beta-1}}\right) = o(n)$. We already observed the behavior for the latter case experimentally in previous works [21, 22]. Thus, an equivalent of the satisfiability threshold conjecture holds for power-law random 2-SAT with power-law exponents $\beta \geq 3$.

6.3 Geometric Distribution

Ansótegui et al. [3] also considered an ensemble of geometric distributions with

$$p_i = \frac{b \cdot (1 - b^{-1/n})}{b - 1} \cdot b^{-(i-1)/n}$$

for $i = 1, \dots, n$ and for some constant $b > 1$. Again, it already holds that $p_1 \geq p_2 \geq \dots \geq p_n$. It holds that

$$p_1^2 = \frac{b^2 \cdot (1 - b^{-1/n})^2}{(b - 1)^2}$$

and

$$\sum_{i=1}^n p_i^2 = \frac{b+1}{b-1} \cdot \frac{1 - b^{-1/n}}{1 + b^{-1/n}}.$$

One can show that $p_1^2 = o(\sum_{i=1}^n p_i^2)$. Theorem 1.1 now tells us that there is a sharp threshold at $\frac{b-1}{b+1} \cdot \frac{1+b^{-1/n}}{1-b^{-1/n}}$. This function grows as fast as $\frac{2 \cdot (b-1)}{(b+1) \cdot \ln b} \cdot n$ in the limit. Thus, an equivalent of the satisfiability threshold conjecture also holds for geometric random 2-SAT with $b > 1$.

7 Discussion and Future Work

We showed a dichotomy of coarse and sharp thresholds for the non-uniform random 2-SAT model depending on the variable probability distribution. In the case of a coarse threshold, the coarseness either stems from two variables being present in too many clauses and forming an unsatisfiable sub-formula of size 4 with constant probability or from a snake with three variables which emerges with constant probability. Furthermore we determined the exact position of the satisfiability threshold in the case of a sharp threshold. Hence, our result generalizes the seminal works by Chvatal and Reed [11] and by Goerdt [24] to arbitrary variable probability distributions. It allows us to prove or disprove an equivalent of the satisfiability threshold conjecture for non-uniform random 2-SAT. For example for power-law random 2-SAT, an equivalent of the conjecture holds for power law exponents $\beta \geq 3$ and the satisfiability threshold is at exactly $\frac{(\beta-3) \cdot (\beta-1)}{(\beta-2)^2} \cdot n$ for $\beta > 3$ and exactly at $4 \cdot \frac{n}{\ln n}$ for $\beta = 3$.

The grand goal of our works is to show similar results for higher values of k , where we already made a first step by showing sharpness for certain variable probability distributions [20]. Another direction we are interested in for $k \geq 3$ is proving bounds on the average computational hardness of formulas around the threshold, for example by showing resolution lower bounds like Mull et al. [38].

References

- [1] D. Achlioptas, L. M. Kirousis, E. Kranakis, and D. Krizanc. Rigorous results for random (2+p)-sat. *Theor. Comput. Sci.*, 265(1-2):109–129, 2001.
- [2] C. Ansótegui, M. L. Bonet, and J. Levy. On the structure of industrial SAT instances. In *15th Intl. Conf. Principles and Practice of Constraint Programming (CP)*, pages 127–141, 2009.
- [3] C. Ansótegui, M. L. Bonet, and J. Levy. Towards industrial-like random SAT instances. In *21st Intl. Joint Conf. Artificial Intelligence (IJCAI)*, pages 387–392, 2009.
- [4] C. Ansótegui, J. Giráldez-Cru, and J. Levy. The community structure of SAT formulas. In *15th Intl. Conf. Theory and Applications of Satisfiability Testing (SAT)*, pages 410–423, 2012.
- [5] C. Ansótegui, M. L. Bonet, J. Giráldez-Cru, and J. Levy. The fractal dimension of SAT formulas. In *7th Intl. Joint Conf. Automated Reasoning (IJCAR)*, pages 107–121, 2014.
- [6] C. Ansótegui, M. L. Bonet, J. Giráldez-Cru, and J. Levy. On the classification of industrial SAT families. In *18th Intl. Conf. of the Catalan Association for Artificial Intelligence (CCIA)*, pages 163–172, 2015.
- [7] V. Bapst and A. Coja-Oghlan. The condensation phase transition in the regular k-sat model. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2016*, pages 22:1–22:18, 2016.
- [8] Y. Boufkhad, O. Dubois, Y. Interian, and B. Selman. Regular random k -sat: Properties of balanced formulas. *J. Autom. Reasoning*, 35(1-3):181–200, 2005.

- [9] M. Bradonjic and W. Perkins. On sharp thresholds in random geometric graphs. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2014*, pages 500–514, 2014.
- [10] K. Bringmann. Why walking the dog takes time: Frechet distance has no strongly sub-quadratic algorithms unless SETH fails. In *55th Symp. Foundations of Computer Science (FOCS)*, pages 661–670, 2014.
- [11] V. Chvatal and B. Reed. Mick gets some (the odds are on his side). In *33rd Symp. Foundations of Computer Science (FOCS)*, pages 620–627, 1992.
- [12] A. Coja-Oghlan. The asymptotic k -SAT threshold. In *46th Symp. Theory of Computing (STOC)*, pages 804–813, 2014.
- [13] A. Coja-Oghlan and K. Panagiotou. The asymptotic k -SAT threshold. *Advances in Mathematics*, 288:985–1068, 2016.
- [14] A. Coja-Oghlan and N. Wormald. The number of satisfying assignments of random regular k -sat formulas. *CoRR*, abs/1611.03236, 2016.
- [15] S. A. Cook. The complexity of theorem-proving procedures. In *3rd Symp. Theory of Computing (STOC)*, pages 151–158, 1971.
- [16] C. Cooper, A. Frieze, and G. Sorkin. Random 2SAT with prescribed literal degrees. *Algorithmica*, 48:249–265, 2007.
- [17] M. Cygan, J. Nederlof, M. Pilipczuk, M. Pilipczuk, J. M. M. van Rooij, and J. O. Wojtaszczyk. Solving connectivity problems parameterized by treewidth in single exponential time. In *52nd Symp. Foundations of Computer Science (FOCS)*, pages 150–159, 2011.
- [18] J. Díaz, L. M. Kirousis, D. Mitsche, and X. Pérez-Giménez. On the satisfiability threshold of formulas with three literals per clause. *Theoretical Computer Science*, 410(30-32):2920–2934, 2009.
- [19] J. Ding, A. Sly, and N. Sun. Proof of the satisfiability conjecture for large k . In *47th Symp. Theory of Computing (STOC)*, pages 59–68, 2015.
- [20] T. Friedrich and R. Rothenberger. Sharpness of the satisfiability threshold for non-uniform random k -sat. In *21st Intl. Conf. Theory and Applications of Satisfiability Testing (SAT)*, pages 273–291, 2018.
- [21] T. Friedrich, A. Krohmer, R. Rothenberger, T. Sauerwald, and A. M. Sutton. Bounds on the satisfiability threshold for power law distributed random SAT. In *25th European Symposium on Algorithms (ESA)*, pages 37:1–37:15, 2017.
- [22] T. Friedrich, A. Krohmer, R. Rothenberger, and A. M. Sutton. Phase transitions for scale-free SAT formulas. In *31st Conf. Artificial Intelligence (AAAI)*, pages 3893–3899, 2017.
- [23] J. Giráldez-Cru and J. Levy. A modularity-based random SAT instances generator. In *24th Intl. Joint Conf. Artificial Intelligence (IJCAI)*, pages 1952–1958, 2015.
- [24] A. Goerdt. A threshold for unsatisfiability. *J. Comput. Syst. Sci.*, 53(3):469–486, 1996.
- [25] M. T. Hajiaghayi and G. B. Sorkin. The satisfiability threshold of random 3-SAT is at least 3.52. Technical Report RC22942, IBM, October 2003.
- [26] R. Impagliazzo and R. Paturi. On the complexity of k -SAT. *J. Comput. Syst. Sci.*, 62(2):367–375, 2001.
- [27] R. Impagliazzo, R. Paturi, and F. Zane. Which problems have strongly exponential complexity? In *39th Symp. Foundations of Computer Science (FOCS)*, pages 653–663, 1998.
- [28] A. C. Kaporis, L. M. Kirousis, and E. G. Lalas. The probabilistic analysis of a greedy satisfiability algorithm. *Random Struct. Algorithms*, 28(4):444–480, 2006.

- [29] R. M. Karp. Reducibility among combinatorial problems. In *Proceedings of a symposium on the Complexity of Computer Computations, held March 20-22, 1972, at the IBM Thomas J. Watson Research Center, Yorktown Heights, New York.*, pages 85–103, 1972.
- [30] L. A. Levin. Universal sorting problems. *Problems of Information Transmission*, 9:265–266, 1973.
- [31] J. Levy. Percolation and phase transition in SAT. *CoRR*, abs/1708.06805, 2017. URL <http://arxiv.org/abs/1708.06805>.
- [32] M. Mézard, G. Parisi, and R. Zecchina. Analytic and algorithmic solution of random satisfiability problems. *Science*, 297(5582):812–815, 2002.
- [33] D. G. Mitchell, B. Selman, and H. J. Levesque. Hard and easy distributions of SAT problems. In *10th Conf. Artificial Intelligence (AAAI)*, pages 459–465, 1992.
- [34] R. Monasson and R. Zecchina. Statistical mechanics of the random k -satisfiability model. *Phys. Rev. E*, 56:1357–1370, Aug 1997.
- [35] R. Monasson, R. Zecchina, S. Kirkpatrick, B. Selman, and L. Troyansky. Phase transition and search cost in the $2+p$ -sat problem. *4th Workshop on Physics and Computation, Boston, MA, 1996.*, 1996.
- [36] 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 Struct. Algorithms*, 15(3-4):414–435, 1999.
- [37] R. Motwani and P. Raghavan. *Randomized algorithms*. Cambridge university press, 1995.
- [38] N. Mull, D. J. Fremont, and S. A. Seshia. On the hardness of SAT with community structure. In *19th Intl. Conf. Theory and Applications of Satisfiability Testing (SAT)*, pages 141–159, 2016.
- [39] V. Rathi, E. Aurell, L. K. Rasmussen, and M. Skoglund. Bounds on threshold of regular random k -sat. In *13th Intl. Conf. Theory and Applications of Satisfiability Testing (SAT)*, pages 264–277, 2010.