# PHASE TRANSITIONS IN RANDOM $k$-SAT

PRATIK CHAUDHARI*

*Abstract:* This project discusses bounds for sharp transition thresholds in random satisfiability problems using the second moment method. We look at bounds for random $k - $SAT formulas and motivate the structure of the solution space using these arguments. Next, we briefly discuss the phenomenon of dynamical phase transition, i.e., transition into a glassy phase wherein, the solution space has exponentially many small clusters. The glassy phase has been a major bottleneck for local algorithms using message passing to compute satisfying assignments for random $k$-SAT.

## 1. INTRODUCTION

It has become increasingly apparent that phenomenon in statistical physics such as phase transitions and glassy phases have a strong bearing on important problems in computer science and information theory such as satisfiability, error correcting codes etc. Realizing that extreme slowdown in algorithms for these problems and glassy phase transitions occur hand in hand is interesting both theoretically, and also promises to provide insights into faster, more efficient algorithms. This has spurred interest in ideas from statistical physics such as "replica symmetry" and "cavity method", which are the basis of state-of-the-art algorithms like survey propagation. Let us note that these concepts have also triggered development in the compressed sensing literature, e.g, [KMS+12].

Motivated by these results, we look at phase transitions in random $k$-SAT problems. We will use the second moment method to obtain bounds for clause densities at which $k$-SAT becomes unsatisfiable with high probability. This is however not very illuminative, it was noticed that a number of (randomized) polynomial time algorithms do not work close to the satisfiability threshold. Towards this end, we will also discuss dynamical phase transition, which is when the solution space of $k$-SAT changes drastically — it goes from looking like a giant ball to the "error correcting code" regime. We will discuss major ideas driving this work and sketch proofs for these results. Let us first introduce the problem setup for random $k$-SAT.

*Date*: May 3, 2014.
* Laboratory of Information and Decision Systems, MIT.
Email: pratik.ac@gmail.com.

## 2. SETUP

A $k$-clause is a disjunction of $k$ Boolean variables. Given $n$ Boolean variables, a random $k$-CNF formula, $F_k(n, m)$ is formed by conjunction of $m = rn$ such $k$-clauses, selected uniformly and independently with replacement from the set of all $k$-clauses on these $n$ variables. We assume that all formulas are in conditional normal form (CNF) and call this model "random $k$-SAT". Let SAT be the set of all satisfiable formulas. Also, let $r_k = \sup\{r : F_k(n, rm) \in \text{SAT}\}$ and $r_k^* = \inf\{r : F_k(n, rn) \notin \text{SAT}\}$. Note that we always talk of $F_k(n, rn) \in \text{SAT}$ with high probability, i.e., $\lim_{n \to \infty} P(F_k(n, rn) \in \text{SAT}) = 1$ if $r < r_k$.

Note that $r_k < r_k^*$ and in Sec. 3 we will show that, in fact, $r_k = r_k^*(1 - o(1))$. We will also show that for $k > 22$,

$$r_k \geq 2^k \log 2 - (k+1)\frac{\log 2}{2} - O(1). \tag{1}$$

## 3. SATISFIABILITY THRESHOLDS

3.1. **2-SAT.** 2-SAT has a special structure, observe that $\overline{x_1} \vee x_2$ is equivalent to $x_1 \implies x_2$. Given $F_2(n, rn)$, construct a graph with vertices $x_1, \overline{x_1}, x_2, \ldots, \overline{x_n}$ and draw two edges, $(x_1, x_2)$ and $(\overline{x_2}, \overline{x_1})$ for every clause of the form $\overline{x_1} \vee x_2$. It can be shown that $F_2(n, rn) \notin \text{SAT}$ iff there exists a path from some $x_k$ to $\overline{x_k}$ and from $\overline{x_k}$ to $x_k$, i.e., a bicycle. Indeed, if $\sigma$ satisfies $x_k$, then by following the implications in this bicycle, we get that $\sigma$ satisfies $\overline{x_k}$ as well and vice-versa. Also, if there does not exist a bicycle and the formula is not SAT, we can simply set one of the variables to 1 and follow the implications; we can keep doing this because there are no bicycles, and we arrive at a solution of $F_2(n, rn)$. So the assumption that $F_2(n, rn) \in \text{SAT}$ has to be false.

Let the length of this bicycle be $s$ if some path

$$(u, w_1, w_2, \ldots, w_s, v)$$

with $u, v \in \{w_1, \ldots, w_s, \overline{w_1}, \ldots, \overline{w_s}\}$ exists. The probability of this happening, by Markov's inequality, is

$$P(F_2(n, rn) \notin \text{SAT}) \leq \sum_{s=2}^{n} \binom{n}{s} s^2 (2s)^2 \binom{m}{s+1} \left(\frac{1}{4\binom{n}{2}}\right)^{s+1}.$$

By direct summation, it is easy to see that this is $O(1/n)$ iff $r < 1$, i.e., $r_2^* < 1$. This was an example of the first moment method. Using the second moment method, similar to Sec. 3.3, we can get $r_2 > 1$, i.e., the phase transition threshold is simply $r = 1$. We do this by counting the number of snakes, i.e., cycles where $u = \overline{v}$ and $u \in \{w_1, \ldots, w_s\}$. Note the nature of this analysis, to make the second moment work, we have to resort to considering only a subset of all bicycles.

3.2. **Vanilla second moment method.** To get a ball-park estimate, let us compute a simple upper bound for the

phase transition. The number of all possible $k$-clauses is $C_k = 2^k \binom{n}{k}$ while the number of clauses satisfied by a given random assignment is $S_k = (2^k - 1)\binom{n}{k}$. Thus the probability of $F \in$ SAT is $\binom{S_k}{m}/\binom{C_k}{m} < (1 - 2^{-k})^m$, which means that the expected number of satisfying assignments is $2^n(1 - 2^{-k})^{rn} = o(1)$ for $r > 2^k \log 2$ which implies $r_k^* < 2^k \log 2$.

The idea behind the rest of this section is to get a non-zero lower bound in the limit for $P(X > 0)$, where $X$ is the number of solutions of random $k$-SAT. We can then use Chebyshev's inequality, $P(X > 0) \geq E[X]^2/E[X^2]$. It is however instructive to note that we cannot weigh each satisfying assignment equally, the bounds are too loose that way. Let us review why.

For a random formula $F$, let $S(F)$ be the set of satisfying assignments and let $X = |S(F)|$. If $F = c_1 \wedge c_2 \ldots c_m$, we have

$$E[X^2] = E\left[\left(\sum_\sigma \mathbf{1}_{\sigma \in S(F)}\right)^2\right] = E\left[\sum_{\sigma,\tau} \mathbf{1}_{\sigma,\tau \in S(F)}\right]$$

$$= \sum_{\sigma,\tau} E\left[\prod_c \mathbf{1}_{\sigma,\tau \in S(c)}\right] = \sum_{\sigma,\tau}\prod_c E[\mathbf{1}_{\sigma,\tau \in S(c)}]. \quad (2)$$

We can now see that

$$P(\sigma, \tau \in S(c)) = 1 - 2^{1-k} - 2^{-k}\alpha^k = f_s(\alpha) \quad (3)$$

where $\sigma$ and $\tau$ assign the same value to $z = \alpha n$ variables. This follows because if $\sigma \not\models c$, the only way for $\tau \notin c$ is for all $k$ variables in $c$ to lie in the overlap. $f_S$ thus quantifies the correlation between $\sigma, \tau$ being both satisfying assignments of $F$. We now have

$$E[X^2] = 2^n \sum_{z=0}^n \binom{n}{z} f_S(\alpha)^m,$$

whereupon, using $\binom{n}{z} = \left(\alpha^\alpha(1-\alpha)^{1-\alpha}\right)^{-n} poly(n)$, we have

$$E[X^2] \leq 2^n \left(\max_{0 \leq \alpha \leq 1}\left[\frac{f_S(\alpha)^r}{\alpha^\alpha(1-\alpha)^{1-\alpha}}\right]\right)^n poly(n)$$

$$:= \left(\max_{0 \leq \alpha \leq 1} \Lambda_S(\alpha)\right)^n poly(n). \quad (4)$$

At the same time, note that

$$E[X]^2 = \left(2^n\left(1 - 2^{-k}\right)^m\right)^2 = \Lambda_S(1/2)^n,$$

and hence if we have some $\alpha \in [0,1]$ with $\Lambda_S(\alpha) > \Lambda_S(1/2)$ the the second moment is exponentially greater than the square of the first moment and hence we have an exponentially smaller lower bound on $P(X > 0)$. In other words, we need to reduce the contribution to $E[X^2]$ from $\sigma, \tau$ with overlap of $\alpha > 1/2$; recollect that we used a similar "snake" counting method as opposed to just bicycles. Work on Not All Equal-SAT, i.e., every clause has at least one satisfied and one unsatisfied literal gives

some clues. For NAE-SAT, we have

$$P(\sigma, \tau \in \text{NAE-SAT}) = 1 - 2^{2-k} + 2^{1-k}\left(\alpha^k + (1-\alpha)^k\right)$$

which is symmetric about $\alpha = 1/2$ and as a result, $\Lambda_N(\alpha)$ has a maxima at $1/2$ for $r \leq 2^{k-1}\log 2 - 1$. It was shown in [AM02] that for $r \geq 2^{k-1}\log 2$, NAE-SAT is not satisfiable, i.e., the second moment method actually gives bounds within additive constant. We will use this as a motivation to design a weighing function on the set $S(F)$. Roughly, we reduce the contribution to $E[X^2]$ due to the correlation terms from Eqn. (3).

3.3. **Weighted second moment method.** Let $w(\sigma, F)$ be some weighing function and $w(\sigma, F) = 0$ if $\sigma \notin S(F)$. We assume in this section that $w(\sigma, F)$ factors exactly over the clauses and define $X$ as follows.

$$X = \sum_\sigma w(\sigma, F) = \sum_\sigma \prod_c w(\sigma, c).$$

This choice still works for us because if we can prove $E[X^2] = O\left(E[X]^2\right)$, we will have

$$P(|S(F)| > 0) = P(X > 0) \geq \frac{E[X^2]}{E[X]^2} = O(1).$$

Again,

$$E[X]^2 = 2^n \left(E[w(\sigma, c)]\right)^m$$

$$E[X^2] = \sum_{\sigma,\tau}\prod_c E[w(\sigma, c)w(\tau, c)]$$

$$= \sum_{\sigma,\tau}\left(E[w(\sigma, c)w(\tau, c)]\right)^m$$

where the clause $c = \ell_1 \vee \ldots \vee \ell_k$ is random. We want $w(\sigma, c)$ to be independent of variable labels, so let it be $w(\sigma, c) = w(v)$ where $v_i = 1$ if $\sigma \models \ell_i$ and $-1$ otherwise.

If $v \in A = \{-1, 1\}^k$, since the literals are drawn uniformly and independently, we have

$$E[w(\sigma, c)] = \sum_{v \in A} w(v)2^{-k}.$$

Similarly for $\sigma, \tau$ with an overlap of $z = \alpha n$,

$$E[w(\sigma, c)w(\tau, c)] = \sum_{u,v \in A} w(u)w(v)2^{-k}\prod_{i=1}^k \alpha^{\mathbf{1}_{u_i = v_i}}(1-\alpha)^{\mathbf{1}_{u_i \neq v_i}}$$

$$:= \sum_{u,v \in A} w(u)w(v)\Phi_{u,v}(\alpha) := f_w(\alpha).$$

The second moment now looks like

$$E[X^2] = 2^n \sum_{z=0}^n \binom{n}{z} f_w(\alpha)^m$$

$$\leq 2^n \left(\max_{0 \leq \alpha \leq 1}\left[\frac{f_w(\alpha)^r}{\alpha^\alpha(1-\alpha)^{1-\alpha}}\right]\right)^n poly(n)$$

$$:= \left(\max_{0 \leq \alpha \leq 1} \Lambda_w(\alpha)\right)^n poly(n), \quad (5)$$

with $E[X]^2 = \Lambda_w(1/2)^n$. There are two aspects of the above expression which will help us get $E[X^2]/E[X]^2 = O(1)$. $\Lambda_w(\alpha)$ has a global maximum at $1/2$ and the $poly(n)$

factor above is in fact, just $O(1)$. The former requires that $f'_w(1/2) = 0$. We thus have

$$f'_w(\alpha) = \sum_{u,v \in A} w(u)w(v)\Phi_{u,v}(\alpha) \left[\log \Phi_{u,v}(\alpha)\right]'$$

$$= \sum_{u,v \in A} w(u)w(v)\Phi_{u,v}(\alpha) \sum_{i=1}^{k} \left(\frac{\mathbf{1}_{u_i=v_i}}{\alpha} - \frac{\mathbf{1}_{u_i\neq v_i}}{1-\alpha}\right)$$

$$2^{2k-1} f'_w(\alpha) = \sum_{u,v \in A} w(u)w(v)uv$$

$$= \left(\sum_u w(u)u\right)\left(\sum_v w(v)v\right).$$

Therefore for any $w$, we need

$$f'_w(1/2) = 0 \implies \sum_{v \in A} w(v)v = 0. \qquad (6)$$

It is now evident that the real reason why the calculations in Sec. 3.2 failed was because $w_S(\cdot)$ assigns 0 to $(-1,\ldots,-1)$ and $1/(2^k - 1)$ to all other vectors. On the other hand, for NAE-SAT, $w_N(\cdot)$ assigns 0 to both $(-1,\ldots,-1)$ and $(1,\ldots,1)$ and hence satisfies Eqn. (6). Hence, we shall construct such a $w(\cdot)$. In addition to this, to sharpen the results, we need to make sure that $w$ is as close $w_S$ as possible. By a simple summation, Eqn. (6) can also be written as

$$\sum_{v \neq (-1,\ldots,-1)} w(v)(2\,|v| - k) = 0.$$

where $|v|$ is the number of 1s in $v$. We maximize entropy ($w_S$ is uniform and has maximum entropy) subject to this constraint to get $w(v) \propto \lambda^{|v|}$ with $\lambda$ satisfying $(1+\lambda)^{k-1} = 1/(1-\lambda)$.

We therefore work with a weighing function

$$w(\sigma, F) \propto \prod_c \lambda^{L(\sigma,F)} \mathbf{1}_{\sigma \in S(c)} \qquad (7)$$

where $L(\sigma, F)$ is the number of literals in $F$ satisfied by $\sigma$. This results in the following theorem.

**Theorem 1.** *There exists a sequence $\beta_k \to 0$ s.t. $\forall\, k \geq 22$,*

$$r_k \geq 2^k \log 2 - 2(k+1)\log 2 - 1 - \beta_k.$$

**Remark 2.** Note that the second term is off by a factor of 4 as compared to Eqn. (1). This is because of our constraint that $w(\sigma, F)$ factors exactly over all clauses. It is possible to sharpen these bounds as shown in Sec. 3.4. Also, in the following proof, we will neglect improper clauses, i.e., repeated or contradictory literals since their probability is at most $k^2/n$ each, i.e., w.h.p there are at most $o(n)$ of them. Similarly, we also let the model select clauses with replacement since w.h.p there are at most $o(n)$ clauses that contain the same $k$ variables.

*Proof Sketch.* Let $H(\sigma, F)$ be the number of satisfied literal occurrences in $F$ minus the number of unsatisfied literals. For $0 < \gamma \leq 1$, let $X = \sum_\sigma \gamma^{H(\sigma,F)} \mathbf{1}_{\sigma \in S(F)}$. Note

that $H(\sigma, F) = 2L(\sigma, F) - km$ and hence this is the same as Eqn. (7). Fix $\sigma$ and since literals in $c = l_1 \vee \ldots \vee l_k$ are random, we have

$$E\left[\gamma^{H(\sigma,c)} \mathbf{1}_{\sigma \in S(F)}\right] = E\left[\gamma^{H(\sigma,c)}\right] - E\left[\gamma^{-k} \mathbf{1}_{\sigma \notin S(c)}\right]$$

$$:= \psi(\gamma).$$

$$\implies E[X] = (2\psi(\gamma)^r)^n.$$

Similarly, if $\sigma, \tau$ overlap in $\alpha n$ literals, we can get

$$E\left[\gamma^{H(\sigma,l)+H(\tau,l)}\right] = \alpha\left(\gamma^2 + \gamma^{-2}\right)/2 + 1 - \alpha$$

$$E\left[\gamma^{H(\sigma,l)+H(\tau,l)} \mathbf{1}_{\sigma \notin S(c)}\right] = 2^{-k}(\alpha\gamma^{-2} + (1-\alpha))$$

$$E\left[\gamma^{H(\sigma,l)+H(\tau,l)} \mathbf{1}_{\sigma,\tau \notin S(c)}\right] = 2^{-k}(\alpha\gamma^{-2})$$

$$\implies E\left[\gamma^{H(\sigma,c)+H(\tau,c)} \mathbf{1}_{\sigma,\tau \in S(c)}\right] = \frac{f(\alpha)}{2^k(1-\epsilon)^k}. \qquad (8)$$

for some explicit $f(\alpha)$ and $\epsilon = 1 - \gamma^2$. Using the same calculations as in Eqn. (2), we now have

$$E[X^2] = 2^n \sum_{z=0}^{n} \binom{n}{z} \left(\frac{f(\alpha)}{2^k(1-\epsilon)^k}\right)^n. \qquad (9)$$

The right hand side can be bounded by noting that only $\Theta(n^{1/2})$ of the binomial coefficients dominate, i.e., if $g = \frac{f(\alpha)}{\alpha^\alpha + (1-\alpha)^{1-\alpha}}$, and there exists $\alpha_{\max}$ s.t. $g(\alpha_{\max}) > g(\alpha)$ for all $\alpha \in (0,1)$ and $g''(\alpha_{\max}) < 0$ — it can be shown that these conditions are true in our case for $k > 22$

$$r < 2^k \log 2 - 2\log 2(k+1) - 1 - 3/k,$$

— then there exists $C > 0$ s.t.

$$E[X^2] < C\left(\frac{2g(1/2)}{(2-2\epsilon)^{kr}}\right)^n.$$

where $\epsilon = 1 - \gamma^2$. Now observe that

$$E[X]^2 = \left(\frac{2g(1/2)}{(2-2\epsilon)^{kr}}\right)^n \qquad (10)$$

and the proof is complete by Chebyshev's inequality. ∎

3.4. **Sharper bounds using measure transform.** This section hunts down the missing factor of 4 in Thm. 1. From Eqn. (4), we see that the dominant contributions to $E[X^2]$ come from pairs where fewer than half of the literals satisfied. Hence we construct $S_+(F) = \{\sigma \in S(F) : H(\sigma, F) > 0\}$ (cf. proof of Thm. 1). In Lem. 3, we will show that $E[X_+]/E[X] \to 1/2$ a.s. where $X_+ = \sum_{\sigma \in S_+(F)} \gamma^{H(\sigma,F)}$. If this is true, consider Eqn. 8 with some $\epsilon = 1 - \theta^2$ with $\theta^2 \geq \gamma^2$, we have

$$E\left[\gamma^{H(\sigma,F)+H(\tau,F)} \mathbf{1}_{\sigma,\tau \in S_+(F)}\right] \leq E\left[\theta^{H(\sigma,F)+H(\tau,F)} \mathbf{1}_{\sigma,\tau \in S_+(F)}\right]$$

$$\leq E\left[\gamma^{H(\sigma,F)+H(\tau,F)} \mathbf{1}_{\sigma,\tau \in S(F)}\right]$$

$$:= f_2(\alpha, \epsilon)^m.$$

where $f_2(\alpha, \epsilon)$ is the same as $f(\alpha)(2 - 2\epsilon)^{-k}$ from Eqn. (8). Since this holds for any $\epsilon \leq 1 - \gamma^2$, we can write for

$\epsilon_0 = 1 - \gamma^2$

$$\mathrm{E}\left[\gamma^{H(\sigma,F)+H(\tau,F)}\mathbf{1}_{\sigma,\tau\in S_+(F)}\right] \leq \left[\inf_{\epsilon\leq\epsilon_0} f_2(\alpha,\epsilon)\right]^m$$

$$\implies \mathrm{E}[X_+^2] \leq 2^n \sum_{z=0}^{n} \binom{n}{z} \left[\inf_{\epsilon\leq\epsilon_0} f_2(\alpha,\epsilon)\right]^m$$

Again truncate the binomial coefficients and note that from Lem. 3 and Eqn. (10), we have

$$5\mathrm{E}[X_+^2] > \mathrm{E}[X]^2 = g(1/2,\epsilon_0)^n.$$

If we can now find a piecewise constant function $\xi$ such that $g(1/2,\epsilon_0) > f_2(\alpha,\xi(\epsilon))/(\alpha^\alpha + (1-\alpha)^{1-\alpha})$ we can again apply Chebyshev's inequality to each of the piecewise constant pieces to get $\mathrm{P}(X_+ > 0)$ with a new bound

$$r_k > 2^k \log 2 - \frac{\log 2}{2}(k+1) - 1 - 50k^3 2^{-k}. \qquad (11)$$

Please see Lem. 9 in [AP04] for a proof of this idea.

**Lemma 3.** *For some $\gamma$, as $n \to \infty$, we have $\frac{E[X_+]}{E[X]} \to 1/2$.*

*Proof.* This is proved by a classic measure tilting argument. Let $\mathbf{P}$ be the probability assigned to a sequence of random literals that form $F$, i.e., $l_1, \ldots, l_{km}$. Tilt this measure to ensure $\mathrm{E}_\gamma[H(\sigma,c)\,\mathbf{1}_{\sigma\in S(c)}] = 0$ by using

$$\mathbf{P}_\gamma[H(\sigma,l) = 1] = \frac{\gamma}{\gamma + \gamma^{-1}} = \frac{2\gamma}{\gamma + \gamma^{-1}}\mathbf{P}[H(\sigma,l) = 1].$$

Because literals are independent, for a clause $c$, we have

$$\mathbf{P}_\gamma(c) = \frac{2^k \gamma^{H(\sigma,c)}\mathbf{P}(c)}{(\gamma + \gamma^{-1})^k}.$$

Define a new measure $\tilde{\mathbf{P}}_\gamma(c) \propto \mathbf{P}_\gamma(c)\,\mathbf{1}_{\sigma\in S(c)}$. A random formula $F$ will be sampled from this measure with probability

$$\mathbf{P}_\gamma(F) \propto \gamma^{H(\sigma,F)}\,\mathbf{P}(F)\,\mathbf{1}_{\sigma\in S(F)}$$

and since we have again have $\tilde{\mathrm{E}}_\gamma[H(\sigma,c)] = 0$, we get

$$\tilde{\mathbf{P}}_\gamma[H(\sigma,F) \geq 0] \to 1/2.$$

Use the definition of $X, X_+$ to prove the claim. ∎

## 4. DYNAMICAL PHASE TRANSITION

In this section, we describe a few results on the geometry of the solution space of random $k$-SAT. The motivation for these results comes from the following — a very simple algorithm can find satisfying assignments for $r = O(2^k/k)$ with high probability [CF86], it simply assigns a random value to a randomly unassigned variable and simplifies the clauses. In fact an algorithm based on belief propagation can solve instances up to $r = \Theta(2^k \frac{\log k}{k})$ using an idea known as decimation [MRTS07]. However, in Sec. 3, we saw that $k$-SAT is satisfiable for far more after that; it has solutions (w.h.p) until $r = 2^k \log 2 - O(k)$. This apparent discrepancy can be resolved with the following theorem towards which we work in this section.

It turns out that the solution space, i.e., $S(F)$ looks like a giant ball until $r = O(2^k/k)$ and shatters into exponentially many clusters with exponentially small number of solutions in each as $r$ increases. We will use the Hamming distance between two satisfying instances $\sigma, \tau$ and say that they are adjacent if $|\sigma - \tau| = 1$. A region is then the union of connected components of $S(F)$.



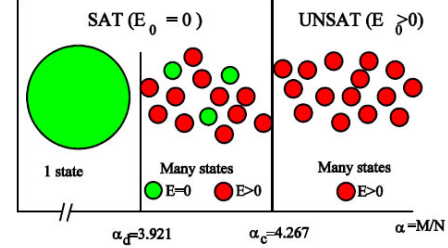FIGURE 1. Phase chart for random 3-SAT

**Theorem 4.** *For $0 < \delta < 1/3$ and $r = (1-\delta)2^k \log 2$, for all $k > k_0(\delta)$, there exist $\alpha < \beta < 1/2$ such that $S(F)$ consists of $2^{\epsilon_k n}$, with $\epsilon_k = \delta/2 - 3/k^2$, non-empty cluster regions with diameter at most $\alpha n$ and distance between every pair of regions at least $\beta n$.*

We will not prove this theorem here, although we shall satisfy ourselves with a rough calculation without proving the details.

4.1. **Exponentially-many clusters.** Note that in Eqn. (4), $\Lambda_S(\alpha) = \mathrm{E}[|S(F)|]$ with overlap ratio of $\alpha$. For convenience, we introduce $\overline{\Lambda} = \Lambda_S(1-\alpha)$, i.e., expected number of pairs of solutions at a Hamming distance of $\alpha n$. The program is as follows: If there exists a $z$ s.t. there are no pairs of assignments at distance $z$, it is an upper bound for the diameter of a cluster. Also, if we can find $\overline{\Lambda} < 1$ in an interval $(\alpha, \beta)$, it is immediate that $S(F)$ can be partitioned into regions (i.e., sets of clusters) with distance at least $\beta n$ and diameter at most $\alpha n$. For a clause $c$ in cluster $C$, let $R(C)$ be assignments with distance at most $\alpha n$ from C and $B(C)$ be all the assignments with distance at most $\beta n$ from $R(C)$. Note that since $\overline{\Lambda} < 1$ for $(\alpha, \beta)$, $B(C) \setminus R(C)$ does not have any satisfying assignments, i.e., the size of the "cluster region" is $\alpha n$ and the distance between any two regions is at least $\beta n$.

We now show that there are exponentially-many such regions. From Sec. 3.3, we have $\mathrm{E}[X^2] < C \max_{\alpha\in[0,1]} \Lambda_w(\alpha)^n$ and since from the second moment method gives $\mathrm{E}[X^2] < C\mathrm{E}[X]^2$, we use the Payley-Zigmund inequality for $t \leq \mathrm{E}[X]$ to get

$$\mathrm{P}(X > t) \geq \frac{(\mathrm{E}[X] - t)^2}{\mathrm{E}[X^2]}.$$

Take $t = \mathrm{E}[X]/poly(n)$ and see that $X$ is within a polynomial factor of its expectation, i.e., $\Lambda_w(1/2)^{n/2}$ with constant probability. It turns out that the event "$F$ has more than $q$ solutions" has a sharp threshold [ACO08] which

implies that for $r < 2^k \log 2 - k$,

$$|S(F)| > \Lambda_w(1/2)^{n/2}/poly(n) \quad \text{w.h.p.}$$

We now divide this by the upper bound of the cluster diameter, let $\Delta = \inf\left\{\alpha : \overline{\Lambda} < 1\right\}$ and

$$g_k = \max_{\alpha \in [0,\Delta]} \overline{\Lambda}(\alpha, k).$$

The expected number of pairs of solutions with distance *at most* $\Delta n$ is then $B < poly(n)\, g_k^n$, since $\overline{\Lambda}$ is expected pairs at particular distance $\alpha n$ and there are most $n + 1$ distances. By Markov's inequality, this means that w.h.p. the number of pairs of solutions at dist. $\Delta n$ is $poly(n)g_k^n$. Since every region has size at most $\Delta n$, w.h.p. the number of pairs in each region is $poly(n)g_k^n$. and if we can show that $g_k < \Lambda_w(1/2)$ we can see that $S(F)$ has at least

$$\frac{1}{poly(n)} \left(\frac{\Lambda_w(1/2)}{g_k}\right)^{n/2}$$

clusters. The rest of the argument then shows that this happens for $\alpha = 1/k$ and $\beta = 1/2 - 5/6\sqrt{\delta}$ in Thm. 4. We do not discuss it here because it is just careful arithmetic, the case for $k > 15$ is analytical and $15 \geq k \geq 8$ is computational. Let us however note that Thm. 4 is quite illuminative, as $r$ approaches the threshold $2^k \log 2$, the clusters are maximally far apart, in fact their size $1/k$ decreases and they vanish as soon as we cross this threshold. Fig. 1 shows this phase transition phenomenon for 3-SAT.

Let us note that, using the results in this section, we can show [ACO08] that there exists a $\epsilon_k \to 0$ s.t. the "glassy" phase with exponentially-many clusters exists for

$$(1 + \epsilon_k)\frac{2^k}{k} \log 2 \leq r \leq (1 - \epsilon_k)\, 2^k \log 2.$$

4.2. **Survey propagation.** As we saw in Sec. 4.1, the solution space of $k$-SAT near the transition threshold lies in the "error correcting code regime" or the "glassy" phase, i.e., it consists of exponentially many clusters with exponentially few solutions in each cluster. Note that, random $k$-SAT is not a hard problem per se, below the dynamical transition threshold, exponentially many solutions exist. It turns out that local algorithms based on belief propagation (BP), warning propagation (WP) etc. work well when there are short-range correlations among variables in the factor graph in SAT. Roughly, BP assumes that the variables in a clause are uncorrelated if that clause is removed (this is true only for trees). This is true only when the solution space looks like a "giant ball".

On the other hand, as shown in [ART06], a large fraction of variables in every cluster region are "frozen", i.e., they take the same value for all the satisfying assignments in the cluster. Moreover, they correlate strongly

with values of frozen variables in neighboring clusters, which is why message passing algorithms on the random factor graph do not work well.

Survey propagation (SP) [BMZ05] then, is an algorithm based on the cavity method from statistical physics. Roughly, it sends a warning message (i.e., $\mu^\alpha(a, i) = 1$) from a factor $a$ to a variable $i$, if all the satisfying assignments of some cluster region $\alpha$ are such that they do not satisfy $a$, in other words, variable $i$ has the complete responsibility of satisfying $a$. Every cluster region is picked with a probability proportional to the number of satisfying assignments in it. SP therefore hinges on explicitly introducing long-range correlations among different cluster regions. This method is also known as 1-step replica symmetry breaking in statistical physics.

## 5. SUMMARY

This project discussed phase transition phenomenon in the solution space of random satisfiability problems. We derived bounds for the satisfiability transition threshold using the first and second moment methods for random 2-SAT and $k$-SAT. We also briefly explored the phenomenon of dynamical phase transition, wherein the solution space changes drastically and bifurcates into exponentially many clusters. This was used to discuss survey propagation that is one of the state-of-the-art solvers for random satisfiability problems.

## REFERENCES

[ACO08] Dimitris Achlioptas and Amin Coja-Oghlan. Algorithmic barriers from phase transitions. In *FOCS*, 2008.

[AM02] Dimitris Achlioptas and Cristopher Moore. The asymptotic order of the random k-SAT threshold. In *FOCS*, 2002.

[AP04] Dimitris Achlioptas and Yuval Peres. The threshold for random k-SAT is $2^k \log 2 - o(k)$. *American Mathematical Society*, 2004.

[ART06] Dimitris Achlioptas and Federico Ricci-Tersenghi. On the solution-space geometry of random constraint satisfaction problems. In *STOC*, 2006.

[BMZ05] Alfredo Braunstein, Marc Mézard, and Riccardo Zecchina. Survey propagation: An algorithm for satisfiability. *Random Structures & Algorithms*, 2005.

[CF86] Ming-Te Chao and John Franco. Probabilistic analysis of two heuristics for the 3-satisfiability problem. *SIAM Journal on Computing*, 1986.

[KMS⁺12] Florent Krzakala, Marc Mézard, François Sausset, YF Sun, and Lenka Zdeborová. Statistical physics based reconstruction in compressed sensing. *Physical Review X*, 2012.

[MRTS07] Andrea Montanari, Federico Ricci-Tersenghi, and Guilhem Semerjian. Solving constraint satisfaction problems through belief propagation-guided decimation. *arXiv:0709.1667*, 2007.