

# Reconstruction and Clustering in Random Constraint Satisfaction Problems

Andrea Montanari\*

Ricardo Restrepo<sup>†</sup>

Prasad Tetali<sup>‡</sup>

December 1, 2009

## Abstract

Random instances of Constraint Satisfaction Problems (CSP's) appear to be hard for all known algorithms, when the number of constraints per variable lies in a certain interval. Contributing to the general understanding of the structure of the solution space of a CSP in the satisfiable regime, we formulate a set of technical conditions on a large family of random CSP's, and prove bounds on three most interesting thresholds for the density of such an ensemble: namely, the *satisfiability* threshold, the threshold for *clustering* of the solution space, and the threshold for an appropriate *reconstruction* problem on the CSP's. The bounds become asymptotically tight as the number of degrees of freedom in each clause diverges. The families are general enough to include commonly studied problems such as, random instances of Not-All-Equal-SAT,  $k$ -XOR formulae, hypergraph 2-coloring, and graph  $k$ -coloring. An important new ingredient is a condition involving the Fourier expansion of clauses, which characterizes the class of problems with a similar threshold structure.

---

\* Department of Electrical Engineering and Department of Statistics, Stanford University; research funded in part by the NSF grants CCF-0743978 and DMS-0806211

<sup>†</sup> School of Mathematics, Georgia Tech, Atlanta, GA 30332-0160

<sup>‡</sup> Schools of Mathematics and Computer Science, Georgia Tech, Atlanta, GA 30332-0160; research funded in part by the NSF grants DMS-0701043 and CCF-0910584

# 1 Introduction

Given a set of  $n$  variables taking values in a finite alphabet, and a collection of  $m$  constraints, each restricting a subset of variables, a Constraint Satisfaction Problem (CSP) requires finding an assignment to the variables that satisfies the given constraints. Important examples include  $k$ -SAT, Not All Equal SAT, graph (vertex) coloring with  $k$  colors etc. Understanding the threshold of satisfiability/unsatisfiability for *random* instances of CSPs, as the number of constraints  $m = m(n)$  varies, has been a challenging task for the past couple of decades, with some notable successes (see e.g., [ANP05]). On the algorithmic side, the challenge of *finding* solutions of a random CSP *close to the threshold of satisfiability* (in the regime where solutions are known to exist) remains widely open. All provably polynomial-time algorithms fail well before the SAT to UNSAT threshold.

The attempt to understand this universal failure led to studying the geometry of the set of solutions of random CSPs [MPZ02, AC08], as well as the emergence of long range correlations among variables in random satisfying assignments [KM+07]. These research directions are motivated by two heuristic explanations of the failure of polynomial algorithms: **(1)** The space of solutions becomes increasingly complicated as the number of constraints increases and is not captured correctly by simple algorithms; **(2)** Typical solutions become increasingly correlated and local algorithms cannot unveil such correlations.

By analyzing a large class of random CSP ensembles, this paper provides strong support to the belief that the above phenomena are *generic*, that they are characterized by *sharp thresholds*, and that the *thresholds for clustering and reconstruction differ at most by a subleading term*, where the notion of ‘subleading’ will be made clearer below.

## 1.1 Related work

Building on a fascinating conjecture on the geometry of the set of solutions, statistical physicists have developed surprisingly efficient message passing algorithms to solve random CSPs. For instance, survey propagation [MPZ02, MZ02] has been shown empirically to find solutions of random 3-SAT extremely close to the SAT-UNSAT transition. In order to understand the success of these heuristics, it has become important to study the thresholds for the emergence of so-called *clustering* of solutions – the emergence of an exponential number of sets (or clusters) of solutions, where solutions within a cluster are closer (in the sense of Hamming distance, say), compared to the intra-cluster distance [MMZ05, AR06, AC08]. Moreover, the fact that solutions within a cluster impose long-range correlations among assignments of variables, motivates one to study the so-called reconstruction problem in the context of random CSP’s. Indeed, non-rigorous statistical mechanics calculations imply that the clustering and reconstruction thresholds coincide [MM06, KM+07].

Further, understanding the threshold for (non)reconstruction is also becoming relevant, if not crucial, to understanding the limit of the Glauber dynamics to sample from the set of solutions of a CSP. Indeed non-reconstructibility was proved in [BK+05] to be a necessary condition for fast mixing, and is expected to be sufficient for a large class of ‘sufficiently random’ problems [GM07].

Reconstruction for models on general *graphical models* -including for instance the case of random proper colorings of the vertices of a graph- was first considered in [BK+05]. The problem amounts to understanding the correlation, as measured e.g., through the mutual information, between the color of a vertex  $v$  and the colors of vertices at distance at least  $t$  from  $v$ . In particular, the problem is said to be ‘unsolvable’ if such a correlation decays to 0 with  $t$ . We refer to Section 3 for a precise definition of the reconstruction problem. For a class of models, including the so-called Ising spin glass, the antiferromagnetic Potts model, and proper  $q$ -colorings of a graph, [GM07] derived a general sufficient condition, under which reconstruction for (sparse) random graphs  $G(n, m)$  with  $m = cn$  edges is possible if and only if it is possible for a Galton-Watson tree with independent Poisson( $2c$ ) degrees for each vertex. Moreover, they also verified that the

condition holds for the Ising spin glass and the antiferromagnetic Potts at non-zero temperature, leaving open the case of proper colorings of graphs, which we settle here.

## 1.2 Summary of contributions

It is against this backdrop that we consider certain general families of CSP's – the first dealing with constraints consisting of  $k$ -tuples of binary variables (as in  $k$ -uniform hypergraph 2-coloring or Not-All-Equal (NAE)  $k$ -sat), while the second dealing with  $q$ -colorings of vertices of graphs (which may be seen as an instance of a CSP with  $q$ -ary variables) – and study three important threshold phenomena. Our chief contribution is as follows.

(a) We formulate an easy-to-check set of assumptions under which a general class of constraint satisfaction problems (including the models mentioned above) can be understood rather precisely in terms of the thresholds for satisfiability, clustering and (non)reconstruction phenomena. In particular we verify that the last two thresholds coincide within the precision of our bounds. (See Theorems 3.2 and 3.3 for precise statements.)

(b) We consider tree ensembles (families of random CSP's whose variable-constraint dependency structure takes the form of a tree), and prove optimal bounds on the threshold for reconstruction on trees. These CSP's consist of binary variables, and the constraints are  $k$ -ary, and the bounds are optimal to first order, as  $k$  goes to infinity.

(c) We verify the sufficient condition of [GM07] for proper colorings of graphs, thus extending the reconstruction result for colorings on trees to the same on *sparse* random graphs.

(d) By way of techniques, we make crucial use of the Fourier expansion of the (binary  $k$ -CSP) constraints, after introducing an assumption on the Fourier expansion, as part of the random ensemble under consideration; this is key to being able to characterize the thresholds precisely.

(e) Finally, as illustrative examples, we mention the specific bounds (on various thresholds) that follow for some standard models, such as the NAE  $k$ -SAT,  $k$ -XOR formulae etc.

The organization of the paper is as follows. In Section 2, we give the formal definitions and assumptions of our models. We state our main results in Section 3. In Section 4, we state and prove the optimal bounds for the tree reconstruction problem. In Section 5, we verify the sufficient condition (from [GM07]) for the specific problem of graph proper  $q$ -coloring, thus proving one of our main results – optimal bounds on the (sparse) random graph reconstruction problem for colorings. In Appendix A, we derive a certain technical second moment bound that is needed to prove our theorem on the satisfiability threshold. In Appendix B, we prove various technical results need to complete the proof of the clustering threshold. In Appendix C, certain sharp threshold results are derived making use of recent results of [AC08, CD09], so that we can extend the high-probability-statements derived in the previous appendices to hold with probability tending to one. Further details on what is proved in these appendices appear in Section 3.3, after the precise statement of our main results.

## 2 Definitions

In this section we define a family of random CSP ensembles: problems with constraints involving  $k$ -tuples of binary variables. We further define  $q$ -ary ensembles as a natural extension of the latter. We finally introduce some analytic definitions that will be necessary in order to present our results.

*Binary  $k$ -CSP ensemble.* Given an integer  $n$ ,  $\alpha \in \mathbb{R}_+$ , and a distribution  $p = \{p(\varphi)\}$  over Boolean functions  $\varphi : \{+1, -1\}^k \rightarrow \{0, 1\}$ ,  $\text{CSP}(n, \alpha, p)$  is the ensemble of random CSP's over  $n$  Boolean variables

$\underline{x} = (x_1, \dots, x_n)$  defined as follows. For each  $a \in \{1, \dots, m = n\alpha\}$ , draw  $k$  indices  $i_a(1), \dots, i_a(k)$  independently and uniformly at random in  $[n]$ , and a function  $\varphi_a$  with distribution  $p(\varphi)$ . An assignment  $\underline{x}$  satisfies the resulting instance if  $\varphi_a(x_{i_a(1)}, \dots, x_{i_a(k)}) = 1$  for each  $a \in [m]$ . A CSP instance can be naturally described by a bipartite graph  $G$  (often referred to in the literature as a ‘factor graph’) including a node for each clause  $a \in [m]$  and for each variable  $i \in [n]$ , and an edge  $(i, a)$  whenever variable  $x_i$  appears in the  $a$ -th clause.

*q-ary ensembles.* A  $q$ -ary ensemble is the natural generalization of a binary ensemble to the case in which variables take  $q$  values. For the sake of simplicity, we restrict our discussion here to the case of pairwise constraints (i.e.  $k = 2$  in the language of the previous paragraph).

Given an integer  $n$ ,  $\alpha \in \mathbb{R}_+$ , and a distribution  $p = \{p(\varphi)\}$  over Boolean functions  $\varphi : [q] \times [q] \rightarrow \{0, 1\}$ ,  $\text{CSP}_q(n, \alpha, p)$  is the collection of random CSP’s over  $q$ -ary variables  $x_i$ , for  $i = 1, 2, \dots, n$ , defined as follows. For each  $a \in \{1, \dots, m = n\alpha\}$ , draw 2 indices  $i_a, j_a$  independently and uniformly at random in  $[n]$ , and a function  $\varphi_a$  with distribution  $p(\varphi)$ . An assignment  $\underline{x} = (x_1, \dots, x_n)$  satisfies the resulting instance, if  $\varphi_a(x_{i_a}, x_{j_a}) = 1$  for each  $a \in [m]$ .

In this paper, by way of illustrating how the results for binary ensembles could be (purportedly) extended to  $q$ -ary ensembles, we will study the  $q$ -coloring model which consists of ensembles with the single clause  $\varphi(x, y) = \mathbb{I}(x \neq y)$ . This model corresponds to proper colorings with  $q$  colors of a random sparse graph with an edge-to-vertex density of  $\alpha > 0$ .

In the rest of this section, we briefly review some well known definitions in discrete Fourier analysis that are useful for stating our results. For general background on this material, the reader may consult any classical textbook on (discrete) Fourier analysis or the lecture notes by Diaconis[Dia88]; for a more breezy introduction and a summary of some key tools one may also find the recent survey [Odo08] useful.

*Functional analysis of clauses.* We denote by  $v_\theta$ , the measure defined over  $\{-1, +1\}^k$  such that

$$v_\theta(x) = \prod_{i=1}^k \left( \frac{1 + x_i \theta}{2} \right) \tag{1}$$

for every  $x \in \{-1, +1\}^k$ . This is just the measure induced by choosing  $k$  independent copies of a random variable that takes values  $\pm 1$  and has expectation  $\theta$ . Notice that when  $\theta = 0$ ,  $v_\theta$  corresponds to the uniform measure over  $\{-1, +1\}^k$ .

The inner product induced by this measure, on the space of real functions defined on  $\{-1, +1\}^k$  is denoted by  $(\cdot, \cdot)_\theta$ , and the corresponding norm by  $\|\cdot\|_\theta$ . If  $\theta = 0$ , we drop the subindex and just use  $(\cdot, \cdot)$  and  $\|\cdot\|$ , respectively. Thus, if  $f, g : \{-1, +1\}^k \rightarrow \mathbb{R}$ , then

$$\begin{aligned} (f, g)_\theta &= \sum_{x \in \{-1, +1\}^k} f(x) g(x) v_\theta(x), & \|f\|_\theta^2 &= \sum_{x \in \{-1, +1\}^k} f^2(x) v_\theta(x), \\ (f, g) &= \frac{1}{2^k} \sum_{x \in \{-1, +1\}^k} f(x) g(x), & \|f\|^2 &= \frac{1}{2^k} \sum_{x \in \{-1, +1\}^k} f^2(x). \end{aligned}$$

We denote the Hilbert space of functions  $\{-1, +1\}^k \rightarrow \mathbb{R}$  under the inner product  $(\cdot, \cdot)$  by  $J_k$ .

*Fourier transform of clauses.* For any  $Q \subseteq [k] \equiv \{1, \dots, k\}$ , let  $\gamma_Q(x) \stackrel{\text{def}}{=} \prod_{i \in Q} x_i$ . Under the scalar product defined above (with  $\theta = 0$ ), the functions  $\{\gamma_S\}_{S \subseteq [k]}$  form an orthonormal basis for  $J_k$ . Moreover,

they are exactly the algebraic characters of  $\{-1, 1\}^k$  with the group operation of pointwise multiplication. Thus, we define the Fourier transform of a function  $f \in J_k$ , by letting for any  $Q \subseteq [k]$ ,

$$f_Q \stackrel{\text{def}}{=} (\gamma_Q, f) = 2^{-k} \sum_{x \in \{-1, +1\}^k} f(x) \gamma_Q(x).$$

*Noise operator.* Given  $\theta \in [-1, 1]$ , we recall the *Bonami-Beckner* operator  $T_\theta : J_k \rightarrow J_k$  [Bon70, Bec75], by

$$(T_\theta f)(x) \stackrel{\text{def}}{=} \sum_{y \in \{-1, 1\}^k} f(xy) v_\theta(y),$$

where  $xy = (x_1 y_1, \dots, x_k y_k)$ . Notice that  $(T_\theta f)(x)$  corresponds to the expected value of  $f(\mathbf{x}_\theta)$ , where  $\mathbf{x}_\theta$  is obtained from  $x$  by flipping each coordinate independently with probability  $(1 - \theta)/2$ . Notice that  $T_1$  is just the identity operator and  $T_0$  sends  $f$  to the constant function  $(f, \gamma_\emptyset)$ .

The Bonami-Beckner operator diagonalizes with respect to the Fourier basis, in the sense that  $(T_\theta \gamma_Q)(x) = \theta^{|Q|} \gamma_Q(x)$  for any  $Q \subseteq [k]$ .

More generally, given  $h \in [-1, 1]^k$ , we define  $(T_h f)(x) \stackrel{\text{def}}{=} \mathbb{E}[f(\mathbf{x}_h)]$ , where  $\mathbf{x}_h$  is obtained from  $x$  by flipping the  $i^{\text{th}}$  coordinate independently and with probability  $\frac{1-h_i}{2}$ . Since  $T_h$  also diagonalizes with respect to the Fourier basis, one gets  $(T_h \gamma_S)(x) = \gamma_S(h) \gamma_S(x)$ .

*Discrete derivative and influence.* Given a function  $f \in J_{k-1}$ , we define its *discrete derivative*  $f^{(1)} \in J_{k-1}$  as  $f^{(1)}(x) = \frac{1}{2} [f(1, x) - f(-1, x)]$ . We define analogously  $f^{(i)}$  for any other variable index. Finally, the *influence* of the  $i^{\text{th}}$  variable on  $f$  is defined using the norm of the derivative

$$I_i(f) \stackrel{\text{def}}{=} \left\| f^{(i)} \right\|^2.$$

For any  $Q \subseteq [k]$ ,  $f_Q^{(i)} = f_{Q \cup \{i\}}$ .

### 3 Main results

As mentioned in the introduction, our goal is in estimating the thresholds for satisfiability, clustering and reconstruction in random CSP's. In general, one should speak of threshold functions depending on the problem size  $n$ . With a slight abuse of notation, we shall leave implicit the dependence on  $n$  of threshold functions unless necessary.

#### 3.1 Binary $k$ -CSP ensembles

##### 3.1.1 Assumptions

We assume the following conditions on the ensemble.

1. *Permutation symmetry.* If  $\varphi^\pi$  is the Boolean function obtained from  $\varphi$  by permuting its arguments, we require  $p(\varphi^\pi) = p(\varphi)$ . (Notice that this assumption does not imply any loss of generality. Indeed, in the definition of the ensemble  $\text{CSP}(n, \alpha, p)$  the indexes of the arguments of clause  $\varphi_a(x_{i_a(1)}, \dots, x_{i_a(k)})$  are independent and uniformly random.

2. *Balance.* The distribution  $p$  is supported on Boolean functions such that  $\varphi(x_1, \dots, x_k) = \varphi(-x_1, \dots, -x_k)$ . This condition implies that the odd Fourier coefficients of  $\varphi$  are zero.

**3. Feasibility.** For each Boolean function  $\varphi$  in the support of  $p$ , every partial assignment  $(x_1, \dots, x_{k-1})$  can be extended to a satisfying assignment  $(x_0, x_1, \dots, x_{k-1})$  of  $\varphi$ . This condition implies that  $\|\varphi\|^2 \geq 1/2$ .

**4. Dominance of balanced assignments.** For every  $\theta \in [-1, 1]$ ,

$$\mathbb{E}_\varphi \log \|\varphi\|_\theta \leq \mathbb{E}_\varphi \log \|\varphi\|,$$

with equality if and only if  $\theta = 0$ . This condition implies that, in a typical random instance, most solutions are balanced in the sense that they have almost as many +1's as -1's.

While our ultimate goal is to exhibit results as  $k \rightarrow \infty$ , the probability distribution  $p$  over the functions  $\varphi : \{-1, 1\}^k \rightarrow \{0, 1\}$  must be defined for *every*  $k$ , and some agreement should exist between such probability distributions for different  $k$ 's. In our work this agreement is given by two conditions concerning the derivative of the clauses in the support of  $p$ :

(a)  $\ell_1$  norm of the Fourier transform grows at most polynomially in  $k$ . That is, for every  $\varphi \in \text{supp}(p)$ ,

$$\sum_Q \left| \varphi_Q^{(i)} \right| \leq k^a, \quad (2)$$

for some constant  $a$  not depending on  $k$ , and recall that  $\varphi_Q^{(i)} = (\gamma_Q, \varphi^{(i)})$ .

(b) 'Small weight' Fourier coefficients are small. There is a constant  $C > 0$  (not depending on  $k$ ) such that for every  $\varphi \in \text{supp}(p)$ ,

$$\left\| \mathbb{T}_\theta \varphi^{(i)} \right\|^2 \leq e^{-Ck(1-\theta)} \left\| \varphi^{(i)} \right\|^2, \quad \theta \in [0, 1]. \quad (3)$$

### 3.1.2 A few remarks

The feasibility conditions implies that all the variables of  $\varphi$  have the same influence, namely,

$$I_i(\varphi) = \frac{1 - \|\varphi\|^2}{2}. \quad (4)$$

In order to prove this consider, say,  $i = 1$  and let  $N_{ab}(\varphi)$   $a, b \in \{0, 1\}$  be the number of partial assignments  $x_1, \dots, x_{k-1}$  such that  $\varphi(+1, x_1, \dots, x_{k-1}) = a$  and  $\varphi(-1, x_1, \dots, x_{k-1}) = b$ . Then, by definition we have

$$\|\varphi\|^2 = \frac{1}{2^k} [N_{01}(\varphi) + N_{10}(\varphi) + 2N_{11}(\varphi)], \quad (5)$$

$$I_1(\varphi) = \frac{1}{2^{k+1}} [N_{01}(\varphi) + N_{10}(\varphi)], \quad (6)$$

whence our claim (4) follows using  $N_{01}(\varphi) + N_{10}(\varphi) + 2N_{11}(\varphi) = 2^{k-1}$ .

Condition (a) above on the  $\ell_1$  norm of the Fourier transform implies in particular, that for any fixed  $l$ , there exists  $A_l > 0$  (independent of  $k$ ), such that

$$\sum_{1 \leq |Q| \leq l} |\varphi_Q|^2 \leq A_l e^{-Ck/2} \sum_{|Q| \geq 1} |\varphi_Q|^2. \quad (7)$$

An equivalent formulation of Eq. (3), with a possibly different constant  $C$ , is

$$\left( \mathbb{T}_\theta \varphi^{(i)}, \varphi^{(i)} \right) \leq e^{-Ck(1-\theta)} \left\| \varphi^{(i)} \right\|^2, \quad \theta \in [0, 1]. \quad (8)$$

### 3.1.3 Results

An ensemble of binary  $k$ -CSP's will be characterized by the following quantities.

$$\frac{1}{\Omega_k} \stackrel{\text{def}}{=} \mathbb{E}_\varphi \frac{2\mathbf{I}_1(\varphi)}{1 - 2\mathbf{I}_1(\varphi)}, \quad \frac{1}{\tilde{\Omega}_k} \stackrel{\text{def}}{=} -\mathbb{E}_\varphi \log\left(1 - 2\mathbf{I}_1(\varphi)\right), \quad \frac{1}{\tilde{\tilde{\Omega}}_k} \stackrel{\text{def}}{=} \frac{2\mathbb{E}_\varphi \mathbf{I}_1(\varphi)}{1 - 2\mathbb{E}_\varphi \mathbf{I}_1(\varphi)}.$$

Notice that  $\Omega_k \leq \hat{\Omega}_k$ , and that  $\Omega_k \leq \tilde{\tilde{\Omega}}_k$ . Indeed, the first inequality follows by using the inequality  $\log(z) \leq z - 1$  with  $z = 1/(1 - 2\mathbf{I}_1)$ , and the second follows by Jensen's, noting the convexity of  $x \mapsto (2x)/(1 - 2x)$ . More over,  $\hat{\Omega}_k \approx (e^{1/\hat{\Omega}_k} - 1)^{-1} \leq \tilde{\tilde{\Omega}}_k$ ; indeed, letting  $X = -\log(1 - 2\mathbf{I}_1(\varphi))$ , and using Jensen's, we have:

$$\frac{1}{\tilde{\tilde{\Omega}}_k} = \frac{\mathbb{E}(1 - e^{-X})}{\mathbb{E}e^{-X}} = \frac{1}{\mathbb{E}e^{-X}} - 1 \leq e^{\mathbb{E}(X)} - 1 = e^{1/\hat{\Omega}_k} - 1.$$

**Proposition 3.1** *A random binary constraint satisfaction instance from the  $\text{CSP}(n, \alpha, p)$  ensemble is satisfiable, with high probability, if  $\alpha < \alpha_s(k)(1 - o_n(1))$ , where*

$$\Omega_k \log 2 \{1 + o_k(1)\} \leq \alpha_s(k, n) \leq \hat{\Omega}_k \log 2 \{1 + o_k(1)\}.$$

*Vice versa, if  $\alpha > \alpha_s(k)(1 + o_n(1))$ , then with high probability, a  $\text{CSP}(n, \alpha, p)$  instance is unsatisfiable. Further  $|\Omega_k^{-1} - \hat{\Omega}_k^{-1}| \leq 8\mathbb{E}_\varphi \{\mathbf{I}_1(\varphi)^2\}$ .*

As clarified by the last part of the statement, the upper and lower bound approach each other when the influence of a single variable in a clause becomes smaller.

Given a measure  $\mu(\underline{x})$  over variable assignments in  $\{+1, -1\}^V$ , the reconstruction problem is said to be unsolvable if correlations with respect to  $\mu$  decay rapidly with the distance  $r$  on  $G$ . More precisely, if  $\mu_{i, \sim r}$  denotes the joint distribution of  $x_i$  and  $\{x_j : d_G(i, j) \geq r\}$ , then  $\lim_{r \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbf{E} \|\mu_{i, \sim r} - \mu_i \mu_{\sim r}\|_{\text{TV}} = 0$ .

**Theorem 3.2** *Let  $\mu(\underline{x})$  be the uniform measure over solutions of an instance from the  $\text{CSP}(n, \alpha, p)$  ensemble. The reconstruction problem is solvable for  $\mu$  if  $\alpha > \alpha_r(k)$ , where*

$$\alpha_r(k) = \frac{\Omega_k}{k} \{\log k + o(\log k)\}.$$

*Vice versa, the reconstruction problem is unsolvable if  $\alpha < \alpha_r(k)$ .*

Given an instance of  $\text{CSP}(n, \alpha, p)$ , a  $d_{\max}$ -cluster of solutions is any equivalence class of solutions under the (closure of the) relation  $\underline{x} \simeq \underline{x}'$  if  $d_{\text{Hamming}}(\underline{x}, \underline{x}') \leq d_{\max}$ . We say that the set of solutions is *clustered* if it is partitioned into exponentially many clusters for some function  $d_{\max} = d_{\max}(n)$  with  $d_{\max}(n) \uparrow \infty$  as  $n \rightarrow \infty$ .

In order to establish clustering, we require two more conditions.

(a') First, a slightly stronger form of *dominance of balanced assignments*:

$$\mathbb{E}_\varphi \{|\varphi|_\theta^2\} \leq \mathbb{E}_\varphi \{|\varphi|^2\}. \quad (9)$$

(b') The following condition on the Fourier transform of clauses

$$\sum_{Q_1 \subseteq Q_2} \mathbb{E}_\varphi \{\varphi_{Q_1} \varphi_{Q_2}\} \theta^{|Q_1|} \delta^{|Q_2| - |Q_1|} \leq \sum_Q \mathbb{E}_\varphi \{\varphi_Q^2\} \theta^{|Q|}, \quad (10)$$

holding for all  $\theta \in [-1, +1]$ ,  $\delta \in [0, 1 - |\theta|]$ . In particular the latter condition holds whenever  $p(\varphi^{(s)}) = p(\varphi)$  for all  $s = (s_1, \dots, s_k) \in \{+1, -1\}^k$ , where  $\varphi^{(s)}(x_1, \dots, x_k) = \varphi(s_1 x_1, \dots, s_k x_k)$ , condition typically known in the literature [CD04], as closure under *polarization*.

**Theorem 3.3** Consider a  $\text{CSP}(n, \alpha, p)$  ensemble satisfying the above conditions, The set of solutions of a random instance from this ensemble is clustered, with high probability, if  $\alpha > \alpha_d(k)$ , where

$$\alpha_d(k) = \frac{\tilde{\Omega}_k}{k} \{\log k + o(\log k)\}.$$

Further  $|\tilde{\Omega}_k^{-1} - \Omega_k^{-1}| \leq 8\mathbb{E}_\varphi\{\mathbb{I}_1(\varphi)^2\}$ .

Thus, a key result of the present paper is that, for a large number of ensembles,  $\alpha_d(k)$  and  $\alpha_r(k)$  differ at most by a quantity whose relative size is negligible for large  $k$ .

**Example 1: 2-coloring hypergraphs.** Let us consider the ensemble of CSP's consisting of clauses of the type  $\varphi$ , where  $\varphi(x_1, \dots, x_k) = \mathbb{I}(\sum x_i \notin \{-k, k\})$ . The  $\text{CSP}(n, \alpha, p)$  in this case, corresponds to the distribution of 2-colorings of a random hypergraph on  $n$  vertices and  $\alpha n$  edges, with edge size  $k$ , and each edge chosen independently and uniformly at random.

The conditions 1-3 (permutation symmetry, balance, feasibility) clearly hold for this model. The dominance of balanced assignments, in its weak and strong form, follows after checking that  $\|\varphi\|_\theta^2 = 1 - \left(\frac{1+\theta}{2}\right)^k - \left(\frac{1-\theta}{2}\right)^k$  is maximized at  $\theta = 0$ . To establish condition **(a)**, cf. Eq. (2), notice that

$$\varphi_Q^{(i)} = -\frac{1}{2^k} [1 - (-1)^{|Q|}],$$

which clearly implies that the  $\ell_1$  norm of the Fourier transform is bounded. In order to check condition **(b)**, cf. Eq. (3), notice that

$$\frac{(\mathbb{T}_\theta \varphi^{(i)}, \varphi^{(i)})}{\|\varphi^{(i)}\|^2} = \left(\frac{1+\theta}{2}\right)^{k-1} - \left(\frac{1-\theta}{2}\right)^{k-1} \leq e^{-k(1-\theta)/2},$$

for all  $\theta \in [0, 1]$ . On the other hand, we have that

$$\begin{aligned} & \sum_{Q_1 \subseteq Q_2} \mathbb{E}_\varphi\{\varphi_{Q_1} \varphi_{Q_2}\} \theta^{|Q_1|} \delta^{|Q_2| - |Q_1|} \\ &= \left(1 - \frac{1}{2^{k-1}}\right) - \frac{1}{2^k} \left[(1+\delta)^k + (1-\delta)^k\right] \\ &+ \left(\frac{1}{2^k}\right)^2 \left[(1+(\delta+\theta))^k + (1-(\delta+\theta))^k + (1+(\delta-\theta))^k + (1-(\delta-\theta))^k\right], \end{aligned}$$

and the previous expression reaches its maximum for  $\delta = 0$ . Thus,

$$\sum_{Q_1 \subseteq Q_2} \mathbb{E}_\varphi\{\varphi_{Q_1} \varphi_{Q_2}\} \theta^{|Q_1|} \delta^{|Q_2| - |Q_1|} \leq \left(1 - \frac{1}{2^{k-2}}\right) + \left(\frac{1}{2^k}\right)^2 \left[\frac{(1+\theta)^k + (1-\theta)^k}{2}\right],$$

and the right hand side of the previous formula is equal to  $\sum_Q \mathbb{E}_\varphi\{\varphi_Q^2\} \theta^{|Q|}$ , proving condition **(b')**.

Now, an easy computation shows that  $\Omega_k = \tilde{\Omega}_k = 2^{k-1} - 1$  and  $\tilde{\Omega}_k^{-1} = -\log(1 - 2^{-k+1})$ , therefore we have:

	Reconstruction - Clustering	Lower bound satisfiability	Upper bound satisfiability
2-coloring	$(2^{k-1}/k) [\log k + o(\log k)]$	$2^{k-1} \log 2 [1 + o(1)]$	$2^{k-1} \log 2 [1 + o(1)]$



**Example 2: Not All Equal  $k$ -SAT.** Let us consider now an ensemble of CSP's consisting of clauses of type  $\{\varphi_s\}_{s \in \{+1, -1\}^k}$ , where  $\varphi_s(x_1, \dots, x_k) = \mathbb{I}(\sum x_i s_i \notin \{-k, k\})$  and  $p(\varphi_s) = 2^{-k}$  for each  $s \in \{+1, -1\}^k$ . In this case, the CSP  $(n, \alpha, p)$  model corresponds to the distribution of NAE  $k$ -SAT instances for a random formula in  $n$  variables, consisting of  $\alpha n$  random clauses, each with  $k$  literals.

For this model, the conditions 1-3 are easily verified. The dominance of balanced assignments in its strong form follows from the fact that

$$\mathbb{E}_s \|\varphi\|_\theta^2 = \mathbb{E}_s \left( 1 - \prod_{i=1}^k \frac{1 + s_i \theta}{2} - \prod_{i=1}^k \frac{1 - s_i \theta}{2} \right) = \mathbb{E}_s \|\varphi\|^2,$$

which for instance implies also the dominance of balanced assignments in its weak form:

$$2\mathbb{E}_s \log \|\varphi\|_\theta \leq \log \mathbb{E}_s \|\varphi\|_\theta^2 = \log \mathbb{E}_s \|\varphi\|^2 = 2\mathbb{E}_s \log \|\varphi\|.$$

On the other hand, the Fourier expansion of  $\varphi_s$  is given by  $\varphi_{s,Q} = -2^{-k}[\gamma_Q(s) + \gamma_Q(-s)]$  (for  $Q \neq \emptyset$ ) and  $\varphi_{s,Q}^{(i)} = -2^{-k}\gamma_Q(s)[1 - (-1)^{|Q|}]$ . In particular  $|\varphi_{s,Q}^{(i)}| = 2^{-k}[1 - (-1)^{|Q|}]$ , so that both Eqs. (2) and (3) hold along the same lines as the previous example, while the condition (b') follows from the closure under *polarization* of this model. Indeed, in this case we get the same values for  $\Omega_k$ ,  $\tilde{\Omega}_k$  and  $\hat{\Omega}_k$ , so that, we have:

	Reconstruction - Clustering	Lower bound satisfiability	Upper bound satisfiability
NAE-SAT	$(2^{k-1}/k) [\log k + o(\log k)]$	$2^{k-1} \log 2 [1 + o(1)]$	$2^{k-1} \log 2 [1 + o(1)]$

**Example 3:  $k$ -XOR formulas.** For an even integer  $k$ , the  $k$ -XOR ensemble ( $k$  even) consists of clauses of type  $\{\varphi_\epsilon\}_{\epsilon \in \{+1, -1\}}$ , where  $\varphi_\epsilon = \frac{1}{2}(\gamma_\emptyset + \epsilon\gamma_{[k]})$ . This set of clauses is endowed with the uniform probability distribution  $p(\varphi_{+1}) = p(\varphi_{-1}) = 1/2$ . In this case, the CSP  $(n, \alpha, p)$  model corresponds to a system of  $\alpha n$  random linear equations in  $\mathbb{Z}_2$ , in which every equation involves  $k$  randomly chosen variables (with replacement) from a total of  $n$  possible variables.

Conditions 1-3 hold for  $k$  even, and the dominance of balanced assignments condition in its weak and strong form, follows from the fact that  $\mathbb{E}_\varphi \|\varphi\|_\theta^2 = \mathbb{E}_\varphi \|\varphi\|^2$ . The condition on Fourier expansion of clauses for this model is straightforward: The Fourier expansion of  $\varphi_\epsilon$  is concentrated at  $\emptyset$  and  $[k]$ , so that the Eq. (2) holds with  $a = 0$  and the Eq. (2) holds with  $C = 1$ . Also, condition (b') follows from the following calculation,

$$\sum_{Q_1 \subseteq Q_2} \mathbb{E}_\varphi \{\varphi_{Q_1} \varphi_{Q_2}\} \theta^{|Q_1|} \delta^{|Q_2| - |Q_1|} = \frac{1}{4} + \frac{1}{4}\theta^k = \sum_Q \mathbb{E}_\varphi \{\varphi_Q^2\} \theta^{|Q|}.$$

In this case, we have that  $\Omega_k = 1$ , while  $\hat{\Omega}_k = 1/\log 2$ . Therefore, we have:

	Reconstruction - Clustering	Lower bound satisfiability	Upper bound satisfiability
XOR-SAT	$\frac{1}{k} [\log k + o(\log k)]$	$\log 2 + o(1)$	$1 + o(1)$

We remark here that, in the case of XOR-SAT, the clustering and satisfiability thresholds can be determined *exactly* by exploiting the underlying group structure [MRZ03, CD+03] (see [MM09] for a discussion of the reconstruction problem in XOR-SAT).

### 3.2 $q$ -ary ensembles: graph coloring

The following results concerning the colorability and clustering of proper colorings were proved by Achlioptas and Naor [AN05] and Achlioptas and Coja-Oghlan [AC08], respectively.

**Theorem 3.4** (*Graph  $q$ -colorability [AN05]*) A random graph with  $n$  vertices and  $n\alpha$  edges is satisfiable with high probability if  $\alpha < \alpha_s(q)$ , where

$$\alpha_s(q) = q [\log q + o_q(1)] .$$

Vice versa, if  $\alpha > \alpha_s(q)(1 + o_q(1))$ , such a graph is with high probability uncolorable.

**Theorem 3.5** (*Clustering of  $q$ -colorings [AC08]*) The set of proper  $q$ -colorings of a random graph with  $n$  vertices and  $n\alpha$  edges is clustered with high probability if  $\alpha > \alpha_d(q)$ , where

$$\alpha_d(q) = \frac{q}{2} [\log q + o(\log q)] .$$

One of our main results is to prove a corresponding reconstruction theorem for this model as follows.

**Theorem 3.6** (*Graph  $q$ -coloring reconstruction*) Let  $\mu(\underline{x})$  be the uniform measure over of proper  $q$ -colorings of random graph with  $n$  vertices and  $n\alpha$  edges. For  $q$  large enough, the reconstruction problem is solvable for  $\mu$  if  $\alpha > \alpha_r(q)$ , where

$$\alpha_r(q) = \frac{q}{2} [\log q + \log \log q + O(1)] .$$

Vice versa, the reconstruction problem is unsolvable, with high probability, if  $\alpha < \alpha_r(q)$ .

### 3.3 General strategy

The results described in the previous section are of three types: bounds on the satisfiability thresholds, cf. Proposition 3.1 and Theorem 3.4; on the clustering threshold, cf. Theorems 3.3 and 3.5; on the reconstruction threshold, cf. Theorems 3.2 and 3.6. The proof strategy is as follows.

*The satisfiability threshold* can be upper bounded using the first moment of the number of solutions, and lower bounded using the second moment method. This technique is by now discussed in detail in [AM02, AN05, ANP05]; we describe its application to the general CSP( $n, \alpha, p$ ) ensemble in Appendix A.

*The clustering threshold* can be upper bounded through an analysis of the recursive ‘whitening’ process introduced in [Pa02], and further studied in [BV04, MMW07]. This process associates to each cluster a single configuration in an extended space. This approach was successfully developed in [AR06]. The improved bounds in Theorems 3.3 and 3.5 can be obtained by approximating the CSP ensemble with an appropriate ‘planted’ ensemble [AC08]. The proof of Theorem 3.3 is presented in Appendix B. As mentioned in the introduction, the following sharp threshold statements (in the sense of Friedgut [F05]) are justified in Appendix C: a precise statement of a recent characterization, by Creignou-Daude [CD09], of the class of binary CSP’s for which the satisfiability property exhibits a sharp threshold phenomenon; and also an analogous proof for the property of having an exponential number of solutions; the latter being needed in completing the proof of the clustering threshold.

*The reconstruction threshold* is characterized via a three-step procedure:

(1) Bound the reconstruction threshold for an appropriate ensemble of (infinite) tree instances, i.e. CSP instances for which the associated factor graph is an infinite Galton-Watson tree. In the case of proper  $q$ -colorings, a sharp characterization was obtained independently by two groups in the past year [BVV07, Sly08]. In Section 4 we prove sharp bounds on tree reconstruction for binary CSPs. The proof amounts to deriving an exact distributional recursion for the so-called belief process, and carefully bounding its asymptotic behavior.

(2) Given two ‘balanced’ solutions  $\underline{x}^{(1)}, \underline{x}^{(2)}$  (a solution is balanced if each possible variable value is taken on the same number of vertices), define their *joint type*  $\nu(x, y)$  as the matrix such that the fraction of

vertices  $i$  with  $x_i^{(1)} = x$  and  $x_i^{(2)} = y$  is equal to  $\nu(x, y)$ . Consider the number  $Z_b(\nu)$  of balanced solution pairs  $\{\underline{x}_1, \underline{x}_2\}$  with joint type  $\nu$ . One has to show that  $\mathbb{E} Z_b(\nu)$  is exponentially dominated by its value at the uniform type  $\bar{\nu}(x, y) = 1/q^2$  (with  $q = 2$  for binary CSPs). More precisely  $\mathbb{E} Z_b(\nu) \doteq \exp\{n\Phi(\nu)\}$  with  $\Phi$  achieving its unique maximum at  $\bar{\nu}$ .

This is also a crucial step in the second moment method. It was accomplished in [AN05] for proper  $q$ -colorings of random graphs. In the case of binary CSPs, we prove this estimate in Appendix A.

**(3)** Prove that the above imply that the set of solutions of a random instance is, with high probability, *roughly spherical*. By this we mean that the joint type  $\nu_{12}$  of two uniformly random solutions  $\underline{x}^{(1)}, \underline{x}^{(2)}$  satisfies  $\|\nu_{12} - \bar{\nu}\|_{\text{TV}} \leq \delta$  with high probability for all  $\delta > 0$ . Notice that this implication requires bounding the expected ratio of  $Z_b(\nu)$  to the total number of solution pairs. We prove that the implication nevertheless holds in Section 5 for  $q$ -colorings. The argument for binary CSP's is completely analogous, and we omit it.

Finally, it was proved in [GM07] that, under such a sphericity condition, graph reconstruction and tree reconstruction are equivalent, which finishes the proof of Theorems 3.2 and 3.6.

Notice that the techniques used for the clustering and reconstruction thresholds are very different. Thus it is a surprising (and arguably deep) phenomenon that they do coincide as far as the present techniques can tell.

## 4 Tree ensembles and tree reconstruction for binary $k$ -CSP ensembles

In this section we define tree ensembles and prove estimates about the corresponding tree reconstruction thresholds.

### 4.1 The tCSP( $\alpha, p$ ) ensemble

The ensemble tCSP( $\alpha, p$ ) is defined by  $\alpha \in \mathbb{R}_+$  and a distribution  $p$  over Boolean functions  $\varphi : \{-1, +1\}^k \rightarrow \{0, 1\}$ . We assume the conditions on the distribution  $p$  introduced in Section 3.1. An (infinite) instance from this ensemble is generated starting by a root variable node  $\emptyset$ , drawing an integer  $\eta \stackrel{\mathcal{D}}{=} \text{Poisson}(k\alpha)$  and connecting  $\emptyset$  to  $\eta$  function nodes  $\{1, \dots, \eta\}$ . Each function node has degree  $k$ , and each of its  $k - 1$  descendants is the root of an independent infinite tree. Finally, each function node  $a$  is associated independently, with a random clause  $\varphi$  drawn according to  $p$ .

A uniform solution for such an instance is sampled by drawing the root value  $\mathbf{x}_\emptyset \in \{-1, +1\}$  uniformly at random. The values of descendants of each variable node  $i$  are then drawn recursively. If the function node  $a$  connects  $i$  to  $i_1, \dots, i_{k-1}$ , then the values  $\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{k-1}}$  are sampled uniformly from those that satisfy the clause associated with  $a$ , that is, such that the quantity  $\varphi(x_i, x_{i_1}, \dots, x_{i_{k-1}})$  is equal to 1.

By the *balance* condition, this procedure can be shown to be equivalent to sampling a solution according to the ‘free boundary Gibbs measure.’ The latter is a distribution over solutions of the entire (infinite) tCSP formula defined by considering the uniform distribution over solutions of the first  $\ell$  generations of the tree, and then letting  $\ell \rightarrow \infty$ .

### 4.2 Reconstruction

Given any fixed tree ensemble  $T$ , let  $\mathbf{x}$  be a random satisfying assignment for  $T$  according to the distribution described previously. We denote by  $\mathbf{x}_\ell$  the value of  $\mathbf{x}$  at the variables at generation  $\ell$ , and in the case that the root degree is 1, we denote by  $\mathbf{x}_{0,1}, \dots, \mathbf{x}_{0,k-1}$ , the values at the variable nodes connected to the unique child of the root. Also, we use  $\eta_0$  for the root degree of  $T$ . If the tree ensemble  $T$  has root degree  $\eta_0 = d$ , we denote by  $T_i$ ,  $i = 1, \dots, d$ , the subtree generated by the root, its  $i^{\text{th}}$  child, and the child's descendants.

If  $\eta_0 = 1$ , we denote by  $T'_i$ ,  $i = 1, \dots, k-1$ , the subtree generated by the  $i^{\text{th}}$  child of the root's child and its descendants.

Finally, because the tree ensemble  $T$  could be random (for instance we denote by  $\mathbf{T}$  a random tCSP  $(\alpha, p)$ ), we will use  $\mathbf{E}$  for expectation respect to  $\mathbf{T}$ , and  $\langle \cdot \rangle_T$  for expectation respect to  $\mathbf{x}$  (given  $\mathbf{T} = T$ ) and  $\mathbb{E}$  for expectation respect to any other independent random variable (adding, if not in context, a subindex to indicate such random variable).

*Reconstruction:* For a fixed tree ensemble  $T$ , let  $\mu_{\emptyset, \ell}$  be the joint distribution of  $(\mathbf{x}_0, \mathbf{x}_\ell)$  and let  $\mu_\emptyset, \mu_\ell$  be the marginal distribution of  $\mathbf{x}_0$  and  $\mathbf{x}_\ell$  respectively. The reconstruction rate for  $T$  is defined as the quantity  $\|\mu_{\emptyset, \ell}(\cdot, \cdot) - \mu_\emptyset(\cdot)\mu_\ell(\cdot)\|_{\text{TV}}$ . We say that the reconstruction problem for  $T$  is *tree-solvable* if

$$\liminf_{\ell \rightarrow \infty} \|\mu_{\emptyset, \ell}(\cdot, \cdot) - \mu_\emptyset(\cdot)\mu_\ell(\cdot)\|_{\text{TV}} > 0.$$

Analogously, if  $\mathbf{T}$  is a random tCSP  $(\alpha, p)$ , we define the reconstruction rate of  $\mathbf{T}$  as

$$\mathbf{E} \|\mu_{\emptyset, \ell}(\cdot, \cdot) - \mu_\emptyset(\cdot)\mu_\ell(\cdot)\|_{\text{TV}},$$

and we say that the reconstruction problem for  $\mathbf{T}$  is *tree-solvable*

$$\liminf_{\ell \rightarrow \infty} \mathbf{E} \|\mu_{\emptyset, \ell}(\cdot, \cdot) - \mu_\emptyset(\cdot)\mu_\ell(\cdot)\|_{\text{TV}} > 0.$$

*Bias, compatibility:* Given a satisfying assignment  $x_\ell$  for the variables at generation  $\ell$ , define the ‘bias’ of the root, restricted to the value of the variables at level  $\ell$ , as

$$h_T(x_\ell) \stackrel{\text{def}}{=} \langle \mathbf{x}_0 | \mathbf{x}_\ell = x_\ell \rangle_T.$$

Throughout the next proofs we will study  $h_T(x_\ell)$ , for  $x_\ell$  random and subject to different kind of distributions. Notice that under the balance condition  $\|\mu_{\emptyset, \ell}(\cdot, \cdot) - \mu_\emptyset(\cdot)\mu_\ell(\cdot)\|_{\text{TV}} = \frac{1}{2} \langle |h_T(\mathbf{x}_\ell)| \rangle_T$ . In fact, it is the case that

$$|h_T(x_\ell)|\mu_\ell(x_\ell) = |\mu_{\emptyset, \ell}(1, x_\ell) - \mu_{\emptyset, \ell}(-1, x_\ell)| = 2 \left| \mu_{\emptyset, \ell}(1, x_\ell) - \frac{1}{2}\mu_\ell(x_\ell) \right|$$

and similarly,

$$|h_T(x_\ell)|\mu_\ell(x_\ell) = 2 \left| \mu_{\emptyset, \ell}(-1, x_\ell) - \frac{1}{2}\mu_\ell(x_\ell) \right|.$$

By the balance condition,  $\mu_\emptyset(1) = \mu_\emptyset(-1) = 1/2$ . Therefore,

$$\begin{aligned} \langle |h_T(\mathbf{x}_\ell)| \rangle_T &= \sum_{x_\ell} (|\mu_{\emptyset, \ell}(1, x_\ell) - \mu_\emptyset(1)\mu_\ell(x_\ell)| + |\mu_{\emptyset, \ell}(-1, x_\ell) - \mu_\emptyset(-1)\mu_\ell(x_\ell)|) \\ &= 2 \|\mu_{\emptyset, \ell}(\cdot, \cdot) - \mu_\emptyset(\cdot)\mu_\ell(\cdot)\|_{\text{TV}}. \end{aligned}$$

Now, let  $D_T(x_\ell) \stackrel{\text{def}}{=} \{x\}$  if  $h_T(x_\ell) = x$ ,  $D_T(x_\ell) \stackrel{\text{def}}{=} \{-1, 1\}$  if  $|h_T(x_\ell)| < 1$ . Observe that  $D_T(x_\ell)$  consists of the values of the root that are compatible with the assignment  $x_\ell$  for the variables at generation  $l$ .

*Domain of clauses:* Given a binary function  $\varphi(x_0, \dots, x_{k-1})$ , define the partial solution sets

$$\begin{aligned} S^+(\varphi) &\stackrel{\text{def}}{=} \{(x_1, \dots, x_{k-1}) : \varphi(1, x_1, \dots, x_{k-1}) = 1\}, \\ S^-(\varphi) &\stackrel{\text{def}}{=} \{(x_1, \dots, x_{k-1}) : \varphi(-1, x_1, \dots, x_{k-1}) = 1\}, \end{aligned}$$

$$\Lambda^+(\varphi) \stackrel{\text{def}}{=} S^+(\varphi) \setminus S^-(\varphi), \quad \Lambda^-(\varphi) \stackrel{\text{def}}{=} S^-(\varphi) \setminus S^+(\varphi)$$

If the clause  $\varphi$  is balanced and feasible, we have that  $|S^+(\varphi)| = |S^-(\varphi)| = 2^{k-1} \|\varphi\|^2$  and  $|\Lambda^+(\varphi)| = |\Lambda^-(\varphi)| = 2^k \mathbf{I}_1(\varphi)$ .

**Theorem 4.1** *The reconstruction problem for the ensemble  $tCSP(\alpha, p)$  is tree-solvable if and only if  $\alpha > \alpha_{\text{tree}}(k)$  where*

$$\alpha_{\text{tree}}(k) = \frac{\Omega_k}{k} \{\log k + o(\log k)\}.$$

**Proof.** *Upper bound:*

Given a tree ensemble  $T$ , the rate of ‘naive reconstruction’ for  $T$  is defined as

$$z_\ell(T) \stackrel{\text{def}}{=} \langle \mathbb{I}[h_T(\mathbf{x}_\ell) = 1] \rangle_T \quad (= \langle \mathbb{I}[h_T(\mathbf{x}_\ell) = -1] \rangle_T \text{ by the balance condition}),$$

which indicates the probability that a random assignment for the variables at generation  $\ell$ , distributed as  $\mathbf{x}_\ell$ , fixes the root to be equal to 1 (or  $-1$ ). It is easy to see that  $\langle |h_T(\mathbf{x}_\ell)| \rangle_T \geq z_\ell(T)$ . Observe also, that for any  $x, y \in \{-1, 1\}$ ,

$$\langle \mathbb{I}[h_T(\mathbf{x}_\ell) = x] | \mathbf{x}_0 = y \rangle_T = 2z_\ell(T) \delta_{x,y}. \quad (11)$$

Thus, our objective is to show that in an appropriate regime of the parameter  $\alpha$ , the quantity  $\mathbf{E}[z_\ell(\mathbf{T})]$  remains bounded away from zero as  $\ell \rightarrow \infty$ , implying tree-solvability of the reconstruction problem in such regime. Indeed, this implies tree-solvability by ‘naive reconstruction’, i.e. by the procedure that assigns to the root any value compatible with the values at generation  $\ell$ . By notational convenience, define

$$z_\ell(\alpha) = 2\mathbf{E}[z_\ell(\mathbf{T})] \quad \text{and} \quad \widehat{z}_\ell(\alpha) = 2\mathbf{E}[z_\ell(\mathbf{T}) | \eta_0 = 1].$$

Now, notice that for a tree ensemble  $T$  with root degree  $\eta_0 = d$ , and any assignment  $x_\ell$  for the variables at generation  $\ell$ ,  $h_T(x_\ell) = 1$  iff  $h_{T_i}(x_\ell \upharpoonright T_i) = 1$  for some  $i = 1, \dots, d$ , so that

$$\begin{aligned} 2z_\ell(T) &= \left\langle 1 - \prod_{i=1}^d (1 - \mathbb{I}[h_{T_i}(\mathbf{x}_\ell \upharpoonright T_i) = 1]) \mid \mathbf{x}_0 = 1 \right\rangle_T \\ &= 1 - \prod_{i=1}^d \left\langle (1 - \mathbb{I}[h_{T_i}(\mathbf{x}_\ell) = 1]) \mid \mathbf{x}_0 = 1 \right\rangle_{T_i} \quad (\text{By the tree Markov property}) \\ &= 1 - \prod_{i=1}^d (1 - 2z_\ell(T_i)). \end{aligned}$$

Therefore, averaging over  $T$ , we get

$$\begin{aligned} z_\ell(\alpha) &= \mathbb{E}_\eta \left[ 1 - \prod_{i=1}^\eta (1 - \widehat{z}_\ell(\alpha)) \right], \quad \eta \sim \text{Poisson}(k\alpha) \\ &= 1 - \exp(-k\alpha \widehat{z}_\ell(\alpha)). \end{aligned}$$

On the other hand, given a tree ensemble  $T$  with root degree  $\eta_0 = 1$  and with the clause  $\varphi$  assigned to the root’s child, we have that for any satisfying assignment  $x_\ell$  for the variables at generation  $\ell$ ,  $h_T(x_\ell) = 1$  iff

$$\prod_{i=1}^{k-1} D_{T'_i} \left( x_{\ell-1}^{(i)} \right) \subseteq \Lambda^+(\varphi), \quad (12)$$

where  $x_{\ell-1}^{(i)}$  is the assignment  $x_\ell \upharpoonright T'_i$  for the variables at generation  $\ell - 1$  in the subtree  $T'_i$ . Observe that (12) holds, in particular, if for some  $a = (a_1, \dots, a_{k-1}) \in \Lambda^+(\varphi)$ ,  $h_{T'_i}(x_{\ell-1}^{(i)}) = a_i$  for  $i = 1, \dots, k - 1$ .

Therefore, if  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_{k-1})$  denotes a random uniform vector from  $S^+(\varphi)$ , we have

$$\begin{aligned} z_\ell(T) &\geq \frac{1}{2} \sum_{a \in \Lambda^+(\varphi)} \left\langle \prod_{i=1}^{k-1} \mathbb{I} \left[ h_{T'_i}(\mathbf{x}_{\ell-1}^{(i)}) = a_i \right] \mid \mathbf{x}_0 = 1 \right\rangle_T \\ &= \frac{1}{2} \sum_{a \in \Lambda^+(\varphi)} \mathbb{E}_{\mathbf{y}} \prod_{i=1}^{k-1} \left\langle \mathbb{I} \left[ h_{T'_i}(\mathbf{x}_{\ell-1}) = a_i \right] \mid \mathbf{x}_0 = y_i \right\rangle_{T'_i} \quad (\text{By the tree Markov property}) \\ &= \frac{1}{2} \frac{|\Lambda^+(\varphi)|}{|S^+(\varphi)|} \prod_{i=1}^{k-1} 2^{z_{\ell-1}(T'_i)} \quad (\text{By Eq. (11)}). \end{aligned}$$

This in turn implies, after averaging over  $T$ , that

$$\widehat{z}_\ell(\alpha) \geq \mathbb{E}_\varphi \left[ \frac{2\mathbf{I}_1(\varphi)}{\|\varphi\|^2} \right] (z_{\ell-1}(\alpha))^{k-1} = \frac{(z_{\ell-1}(\alpha))^{k-1}}{\Omega_k},$$

which leads to the recursion  $z_\ell(\alpha) \geq 1 - \exp\left(-k\alpha(z_{\ell-1}(\alpha))^{k-1}/\Omega_k\right)$ . Now, it is standard to verify that this recursion implies that  $z_\ell(\alpha)$  is, for all  $\ell$ , greater or equal than the maximum of the fixed points of the function  $g(z) = 1 - \exp\left(-k\alpha z^{k-1}/\Omega_k\right)$  in the interval  $[0, 1]$ . The minimum value of  $\alpha$  for which such fixed point is positive is given by

$$\alpha^* = \frac{\Omega_k \left(1 + u \left(1 + \frac{1}{u}\right)^{k-2}\right)}{k(k-1)},$$

where  $u$  is the unique solution of the equation  $u = (k-1) \log(1+u)$ . In particular, asymptotically in  $k$ , we have that  $\alpha^* = \frac{\Omega_k}{k} (\log k + o(\log k))$ , which implies the upper bound for  $\alpha_{\text{tree}}$ .

*Lower bound:*

The matching lower bound on  $\alpha_{\text{tree}}(k)$  requires a more elaborate proof; we first prove three lemmas, before returning to complete the lower bound proof.  $\square$

Given a tree ensemble  $T$ , let  $\mathbf{x}_\ell^+ \stackrel{\mathcal{D}}{=} (\mathbf{x}_\ell \mid \mathbf{x}_0 = 1)$  and  $\mathbf{x}_\ell^- \stackrel{\mathcal{D}}{=} (\mathbf{x}_\ell \mid \mathbf{x}_0 = -1)$ . When the tree ensemble is not clear in the definition of  $\mathbf{x}_\ell^+$  (or  $\mathbf{x}_\ell^-$ ), we add a subindex indicating the tree ensemble from where it is defined. Notice that, if  $\mu^+$  and  $\mu^-$  are the distributions of  $\mathbf{x}_\ell^+$  and  $\mathbf{x}_\ell^-$  respectively, then

$$\frac{d\mu^-}{d\mu^+} = \frac{1 - h_T(x_\ell)}{1 + h_T(x_\ell)}. \quad (13)$$

By the balance condition, it's clear that

$$h_T(\mathbf{x}_\ell^+) \stackrel{\mathcal{D}}{=} -h_T(\mathbf{x}_\ell^-). \quad (14)$$

Also, it is easy to show that  $\langle h_T(\mathbf{x}_\ell^+) \rangle_T = \left\langle [h_T(\mathbf{x}_\ell)]^2 \right\rangle_T$  (and therefore  $[R_\ell(T)]^2 \leq \langle h_T(\mathbf{x}_\ell^+) \rangle_T \leq R_\ell(T)$ ), so that non-reconstructibility for  $T$  is equivalent to the condition  $\lim_{\ell \rightarrow \infty} \langle h_T(\mathbf{x}_\ell^+) \rangle_T = 0$  (see [MP03]). Similarly, if  $\mathbf{T}$  is a random tCSP  $(\alpha, p)$  ensemble, non-reconstructibility for  $\mathbf{T}$ , is equivalent to the condition  $\lim_{\ell \rightarrow \infty} \mathbf{E} \left[ \langle h_{\mathbf{T}}(\mathbf{x}_\ell^+) \rangle_{\mathbf{T}} \right] = 0$ .

**Lemma 4.2 (a)** *Given a tree ensemble  $T$  with root degree  $\eta_0 = d$ , we have*

$$\left[ \frac{1 - h_T(\mathbf{x}_\ell^+)}{1 + h_T(\mathbf{x}_\ell^+)} \right] \stackrel{\mathcal{D}}{=} \prod_{i=1}^d \left[ \frac{1 - h_{l,i}}{1 + h_{l,i}} \right], \quad (15)$$

where  $(h_{l,i})_{i=1}^d$  are independent random variables such that  $h_{l,i} \stackrel{\mathcal{D}}{=} h_{T_i}(\mathbf{x}_l^+)$ .

(b) Given a tree ensemble  $T$  with root degree  $\eta_0 = 1$  and with the clause  $\varphi$  assigned to the unique child of the root, we have that

$$\left[ \frac{1 - h_T(\mathbf{x}_{\ell+1}^+)}{1 + h_T(\mathbf{x}_{\ell+1}^+)} \right] \stackrel{\mathcal{D}}{=} \frac{\mathbf{T}_{h_l} \varphi(-1, \mathbf{s})}{\mathbf{T}_{h_l} \varphi(1, \mathbf{s})}, \quad (16)$$

where  $\mathbf{s} \sim \text{Unif}(S^+(\varphi))$  and  $h_l = (h_{l,i})_{i=1}^{k-1}$  are independent random variables such that  $h_{l,i} \stackrel{\mathcal{D}}{=} h_{T'_i}(\mathbf{x}_l^+)$ .

**Proof.** This recursion follows straightforwardly from the recursive definition of tree formulae. The balance condition on clauses implies

$$\frac{1 - h_T(\mathbf{x}_l^+)}{1 + h_T(\mathbf{x}_l^+)} = \frac{\langle \mathbb{I}[\mathbf{x}_l = \mathbf{x}_l^+] | \mathbf{x}_0 = -1 \rangle_T}{\langle \mathbb{I}[\mathbf{x}_l = \mathbf{x}_l^+] | \mathbf{x}_0 = 1 \rangle_T}.$$

Therefore, if the root degree of  $T$  is  $\eta_0 = d$ , we have by the tree Markov property that

$$\frac{1 - h_T(\mathbf{x}_l^+)}{1 + h_T(\mathbf{x}_l^+)} = \prod_{i=1}^d \frac{\langle \mathbb{I}[\mathbf{x}_l = \mathbf{x}_l^+ \upharpoonright T_i] | \mathbf{x}_0 = -1 \rangle_{T_i}}{\langle \mathbb{I}[\mathbf{x}_l = \mathbf{x}_l^+ \upharpoonright T_i] | \mathbf{x}_0 = 1 \rangle_{T_i}},$$

and the last expression has the same distribution as  $\prod_{i=1}^d \frac{1 - u_{l,i}}{1 + u_{l,i}}$ , due to the fact that  $(\mathbf{x}_l^+ \upharpoonright T_i)_{i=1}^d$  are inde-

pendent random assignments for the variables at generation  $l$  of  $T_i$ , such that  $\mathbf{x}_l^+ \upharpoonright T_i \stackrel{\mathcal{D}}{=} \mathbf{x}_{l,T_i}^+$ . This proves

Eq. (15). Now, if the root degree of  $T$  is  $\eta_0 = 1$ , define  $(\tilde{\mathbf{x}}_{l,i}^+)_{i=1}^{k-1}$  to be independent random assignments

for the variables at generation  $l$  of the subtrees  $T'_i$ , such that  $\tilde{\mathbf{x}}_{l,i}^+ \stackrel{\mathcal{D}}{=} \mathbf{x}_{l,T'_i}^+$ . By the tree Markov property, we

have that  $(\mathbf{x}_{l+1}^+ \upharpoonright T'_i)_{i=1}^{k-1} \stackrel{\mathcal{D}}{=} (\mathbf{s}_i \tilde{\mathbf{x}}_{l,i}^+)_{i=1}^{k-1}$  where  $\mathbf{s} \sim \text{Unif}(S^+(\varphi))$ . Using once more the tree Markov property, we get

$$\begin{aligned} \left[ \frac{1 - h_T(\mathbf{x}_{\ell+1}^+)}{1 + h_T(\mathbf{x}_{\ell+1}^+)} \right] &= \frac{\sum_y \varphi(-1, y) \prod_{i=1}^{k-1} \langle \mathbb{I}[\mathbf{x}_l = \mathbf{s}_i \tilde{\mathbf{x}}_{l,i}^+] | \mathbf{x}_0 = y_i \rangle_{T'_i}}{\sum_y \varphi(-1, y) \prod_{i=1}^{k-1} \langle \mathbb{I}[\mathbf{x}_l = \mathbf{s}_i \tilde{\mathbf{x}}_{l,i}^+] | \mathbf{x}_0 = y_i \rangle_{T'_i}} \\ &= \frac{\mathbf{T}_{h_l} \varphi(-1, \mathbf{s})}{\mathbf{T}_{h_l} \varphi(1, \mathbf{s})}, \end{aligned}$$

which is precisely Eq. (16). □

The first step of the above recursion can be analyzed precisely, in terms of its distribution.

**Lemma 4.3** *If  $\mathbf{T}$  is a random tCSP  $(\alpha, p)$  ensemble, then the random variable  $h_{\mathbf{T}}(\mathbf{x}_1^+)$  takes values in  $\{0, 1\}$  and, if  $\alpha < (1 - \delta)(\Omega_k \log k)/k$ , we have  $\mathbf{E} h_{\mathbf{T}}(\mathbf{x}_1^+) \leq 1 - k^{-1+\delta}$ .*

**Proof.** If  $T$  is a tree ensemble with root degree  $\eta_0 = 1$  and clause  $\varphi$  assigned to the root's child, from Part(b) of Lemma 4.2, we have that  $\frac{1 - h_T(\mathbf{x}_1^+)}{1 + h_T(\mathbf{x}_1^+)} \stackrel{\mathcal{D}}{=} \varphi(-1, \mathbf{s})$  where  $\mathbf{s} \sim \text{Unif}(S^+(\varphi))$ . Recall that  $h_{0,i} \equiv 1$ . Therefore,

it follows that  $h_T(\mathbf{x}_1^+) = 1$  with probability  $\frac{|\Lambda^+(\varphi)|}{|S^+(\varphi)|} = 1/\Omega_k$  and  $h_T(\mathbf{x}_1^+) = 0$  otherwise. Similarly, if  $T$

is a tree ensemble with root degree  $\eta_0 = d$ , it follows from Part(a) of Lemma 4.2 that  $h_T(\mathbf{x}_1^+) = 1$  with probability  $1 - (1 - 1/\Omega_k)^d$  and  $h_T(\mathbf{x}_1^+) = 0$  otherwise. This implies then that  $h_T(\mathbf{x}_1^+)$  has support in  $\{0, 1\}$  and that  $\mathbf{E}h_T(\mathbf{x}_1^+) = 1 - \exp(-k\alpha(1 - 1/\Omega_k))$ . The conclusion follows straightforwardly.  $\square$

For subsequent steps we track the averages,  $h_\ell^{\text{ave}} \stackrel{\text{def}}{=} \mathbf{E} \langle h_T(\mathbf{x}_l^+) \rangle_T$  and  $\widehat{h}_\ell^{\text{ave}} \stackrel{\text{def}}{=} \mathbf{E} [\langle h_T(\mathbf{x}_l^+) \rangle_T | \eta_0 = 1]$ , using the following bounds.

**Lemma 4.4** *For any  $\ell \geq 0$  we have*

$$h_\ell^{\text{ave}} \leq 1 - e^{-2k\alpha\widehat{h}_\ell^{\text{ave}}}, \quad \widehat{h}_{\ell+1}^{\text{ave}} \leq \frac{1}{2} F_k(h_\ell^{\text{ave}}) + \frac{1}{2} R_k(\sqrt{h_\ell^{\text{ave}}}), \quad (17)$$

$$F_k(\theta) \stackrel{\text{def}}{=} 2\mathbb{E}_\varphi \left[ \frac{(\varphi^{(1)}, \mathbf{T}_\theta \varphi^{(1)})}{\|\varphi\|^2} \right], \quad R_k(\theta) \stackrel{\text{def}}{=} 2\mathbb{E}_{\varphi_i} \left[ \frac{2\mathbf{I}_1(\varphi)}{\|\varphi\|^2} \sum_{Q \subseteq [k-1]} |(\varphi^{(1)}, \gamma_Q)| \theta^{\max(|Q|, 2)} \right], \quad (18)$$

Finally, if  $h_\ell$  is supported on non-negative values, then

$$\widehat{h}_\ell^{\text{ave}} \leq F_k(h_\ell^{\text{ave}}). \quad (19)$$

**Proof.** We will say that a random variable  $\mathbf{X} \in [-1, +1]$  is ‘consistent,’ if  $\mathbf{E} f(-\mathbf{X}) = \mathbf{E} \left[ \left( \frac{1-\mathbf{X}}{1+\mathbf{X}} \right) f(\mathbf{X}) \right]$  for every function  $f$  such that the expectation values exist. A useful preliminary remark [MM06] is that the random variable  $h_T(\mathbf{x}_l^+)$  is consistent (no matter the tree ensemble). In fact, this follows directly from the Eqs. (13) and (14) above:

$$\begin{aligned} \mathbf{E} f(-h_T(\mathbf{x}_l^+)) &= \sum_{x_l} f(-h_T(x_l)) \mu^+(x_l) = \sum_{x_l} f(-h_T(x_l)) \frac{1 + h_T(x_l)}{1 - h_T(x_l)} \mu^-(x_l) \\ &= \mathbf{E} \left[ f(-h_T(\mathbf{x}_l^-)) \frac{1 + h_T(\mathbf{x}_l^-)}{1 - h_T(\mathbf{x}_l^-)} \right] = \mathbf{E} \left[ f(h_T(\mathbf{x}_l^+)) \frac{1 - h_T(\mathbf{x}_l^+)}{1 + h_T(\mathbf{x}_l^+)} \right] \end{aligned}$$

A number of properties of consistent random variables can be found in [RU08]. Let us now consider the first inequality. If  $T$  is a tree ensemble with root degree  $\eta_0 = d$ , it is immediate to from Eq. (15), that

$$\left\langle \left( \frac{1 - h_T(\mathbf{x}_l^+)}{1 + h_T(\mathbf{x}_l^+)} \right)^{1/2} \right\rangle_T = \prod_{i=1}^d \left\langle \left( \frac{1 - h_{T_i}(\mathbf{x}_l^+)}{1 + h_{T_i}(\mathbf{x}_l^+)} \right)^{1/2} \right\rangle_{T_i}. \quad (20)$$

It is possible to show that consistency implies  $\mathbf{E}X = \mathbf{E}X^2$  and  $\mathbf{E} \left( \frac{1-X}{1+X} \right)^{1/2} = \mathbf{E} \sqrt{1-X^2}$  (through the test functions  $f(x) = x(1+x)$  and  $f(x) = x(1+x)^{1/2}(1-x)^{-1/2}$ ), we thus have

$$\begin{aligned} \sqrt{1 - \langle h_T(\mathbf{x}_l^+) \rangle_T} &= \sqrt{1 - \langle [h_T(\mathbf{x}_l^+)]^2 \rangle_T} \geq \left\langle \sqrt{1 - [h_T(\mathbf{x}_l^+)]^2} \right\rangle_T \quad (\text{by Jensen's ineq.}) \\ &= \left\langle \left( \frac{1 - h_T(\mathbf{x}_l^+)}{1 + h_T(\mathbf{x}_l^+)} \right)^{1/2} \right\rangle_T = \prod_{i=1}^d \left\langle \left( \frac{1 - h_{T_i}(\mathbf{x}_l^+)}{1 + h_{T_i}(\mathbf{x}_l^+)} \right)^{1/2} \right\rangle_{T_i} \\ &= \prod_{i=1}^d \left\langle \sqrt{1 - [h_{T_i}(\mathbf{x}_l^+)]^2} \right\rangle_{T_i} \geq \prod_{i=1}^d \left( 1 - \langle h_{T_i}(\mathbf{x}_l^+) \rangle_{T_i} \right) \quad (\text{using } \sqrt{x} \geq x, \text{ for } x \in [0, 1]). \end{aligned}$$



This implies in particular, if  $\mathbf{T}$  is a random tCSP  $(\alpha, p)$ ,

$$\sqrt{1 - \mathbf{E} \langle h_{\mathbf{T}}(\mathbf{x}_l^+) \rangle_{\mathbf{T}}} \geq \mathbb{E}_{\eta} \left[ \prod_{i=1}^{\eta} (1 - \mathbf{E} [\langle h_{\mathbf{T}}(\mathbf{x}_l^+) \rangle_{\mathbf{T}} | \eta_0 = 1]) \right], \quad \eta \sim \text{Poisson}(k\alpha),$$

from where the first inequality follows.

Now, from the recursion Eq. (16), we have for a tree ensemble  $T$  with root degree  $\eta_0 = 1$ , and random clause  $\varphi$  assigned to the child of the root,

$$h_T(\mathbf{x}_{l+1}^+) = \frac{2 \mathbf{T}_{h_l} \varphi^{(1)}(\mathbf{s})}{1 + \mathbf{T}_{h_l} \psi(\mathbf{s})}, \quad \psi(s) \stackrel{\text{def}}{=} \varphi(1, s) \varphi(-1, s)$$

or alternatively,

$$h_T(\mathbf{x}_{l+1}^+) = \mathbf{T}_{h_l} \varphi^{(1)}(\mathbf{s}) + \left( \mathbf{T}_{h_l} \varphi^{(1)}(\mathbf{s}) \right) \mathcal{G}_k(h_l, \mathbf{s}), \quad \mathcal{G}_k(h_l, s) \stackrel{\text{def}}{=} \left[ \frac{1 - \mathbf{T}_{h_l} \psi(s)}{1 + \mathbf{T}_{h_l} \psi(s)} \right],$$

where  $\mathbf{s} \sim \text{Unif } S^+(\varphi)$ . Notice that for any antisymmetric function  $f(s)$ , we have that  $\mathbb{E}_{\mathbf{s}} f(\mathbf{s}) = \frac{(\varphi^{(1)}, f)}{\|\varphi\|^2}$ . Therefore, due to the fact that  $\mathbf{T}_{h_l} \varphi^{(1)}(s)$  is antisymmetric and  $\mathcal{G}_k(h_l, s)$  is symmetric (both in  $s$  and  $h_l$ , actually), we have the formulas

$$\langle h_T(\mathbf{x}_{l+1}^+) \rangle_T = \frac{2}{\|\varphi\|^2} \left\langle \left( \varphi^{(1)}, \frac{\mathbf{T}_{h_l} \varphi^{(1)}(\mathbf{s})}{1 + \mathbf{T}_{h_l} \psi(\mathbf{s})} \right) \right\rangle_T \quad (21)$$

and

$$\langle h_T(\mathbf{x}_{l+1}^+) \rangle_T = \left\langle \frac{(\varphi^{(1)}, \mathbf{T}_{h_l} \varphi^{(1)})}{\|\varphi\|^2} \right\rangle_T + \left\langle \frac{(\varphi^{(1)}, (\mathbf{T}_{h_l} \varphi^{(1)}) \mathcal{G}_k(h_l, \cdot))}{\|\varphi\|^2} \right\rangle_T. \quad (22)$$

In the last expression, the first term is equal to  $\frac{(\varphi^{(1)}, \mathbf{T}_{\langle h_l \rangle_T} \varphi^{(1)})}{\|\varphi\|^2}$ , while the second term can be written, using Fourier expansion, as

$$\frac{1}{\|\varphi\|^2} \sum_{\substack{Q \subseteq [k-1] \\ |Q| \text{ odd}}} \left( \varphi^{(1)}, \gamma_Q \mathbb{E}_{h_l} [\gamma_Q(h_l) \mathcal{G}_k(h_l, \cdot)] \right) \left( \varphi^{(1)}, \gamma_Q \right).$$

Using the fact that  $\mathbb{E} |\mathbf{X}| \leq (\mathbb{E} \mathbf{X})^{1/2}$  for consistent random variables, we can bound the terms with  $|Q| \geq 3$  by

$$\frac{|(\varphi^{(1)}, 1)|}{\|\varphi\|^2} \sum_{\substack{Q \subseteq [k-1] \\ |Q| \geq 3 \text{ odd}}} \left| (\varphi^{(1)}, \gamma_Q) \right| \left( \prod_{i \in Q} \langle h_{T_i}(\mathbf{x}_l^+) \rangle_{T_i} \right)^{1/2}.$$

Also, using the fact that for any even function  $f(x)$  with  $0 \leq f(x) \leq 1$  and a consistent random variable  $\mathbf{X}$ , we have

$$|\mathbb{E}[\mathbf{X} f(\mathbf{X})]| = |\mathbb{E}[2\mathbf{X}^2 f(\mathbf{X}) / (1 + \mathbf{X}) \mathbb{I}_{\{\mathbf{X} \geq 0\}}]| \leq |\mathbb{E}[2\mathbf{X}^2 / (1 + \mathbf{X}) \mathbb{I}_{\{\mathbf{X} \geq 0\}}]| = |\mathbb{E}[\mathbf{X}]|,$$

we can bound the terms with  $|Q| = 1$ , by

$$\frac{|(\varphi^{(1)}, 1)|}{\|\varphi\|^2} \sum_{i=1}^{k-1} \left( \varphi^{(1)}, \gamma_{\{i\}} \right) \left| \langle h_{T_i}(\mathbf{x}_l^+) \rangle_{T_i} \right|.$$

Therefore, for a random tCSP  $(\alpha, p)$  with root degree  $\eta_0 = 1$ , we obtain after averaging

$$\widehat{h}_{l+1}^{\text{ave}} \leq \mathbb{E}_\varphi \frac{\left(\varphi^{(1)}, \mathbb{T}_{h_l^{\text{ave}}} \varphi^{(1)}\right)}{\|\varphi\|^2} + \mathbb{E}_\varphi \left[ \frac{2\mathbb{I}_1(\varphi)}{\|\varphi\|^2} \sum_{\substack{Q \subseteq [k-1] \\ |Q| \geq 3 \text{ odd}}} \left| \left(\varphi^{(1)}, \gamma_Q\right) \right| \left(\sqrt{h_l^{\text{ave}}}\right)^{\max\{|Q|, 2\}} \right],$$

which is precisely the second inequality in the Lemma.

Now, suppose that  $h_l$  is supported on non-negative values and let  $A_s = \{h_l : \mathbb{T}_{h_l} \varphi^{(1)}(s) > 0\}$ . Notice that the complement of  $A_s$  is  $-A_s$  (due to the antisymmetry of  $\mathbb{T}_{h_l} \varphi^{(1)}(s)$  respect to  $h_l$ ). Therefore, using the consistency of the random variables  $h_{l,i}$ , from the Eq. (21) we get

$$\begin{aligned} \langle h_T(\mathbf{x}_{l+1}^+) \rangle_T &= \frac{2}{\|\varphi\|^2} \left\langle \left( \varphi^{(1)}, \frac{\mathbb{T}_{h_l} \varphi^{(1)}(\mathbf{s})}{1 + \mathbb{T}_{h_l} \psi(\mathbf{s})} \right) \mathbb{I}(h_l \in A_s) - \left( \varphi^{(1)}, \frac{\mathbb{T}_{-h_l} \varphi^{(1)}(\mathbf{s})}{1 + \mathbb{T}_{-h_l} \psi(\mathbf{s})} \right) \mathbb{I}(-h_l \in A_s) \right\rangle_T \\ &= \frac{2}{\|\varphi\|^2} \left\langle \left( \varphi^{(1)}, \frac{\mathbb{T}_{h_l} \varphi^{(1)}(\mathbf{s})}{1 + \mathbb{T}_{h_l} \psi(\mathbf{s})} \right) \mathbb{I}(h_l \in A_s) \left[ 1 - \prod_{i=1}^{k-1} \frac{1 - h_{l,i}}{1 + h_{l,i}} \right] \right\rangle_T \\ &\leq \frac{2}{\|\varphi\|^2} \left\langle \left( \varphi^{(1)}, \mathbb{T}_{h_l} \varphi^{(1)}(\mathbf{s}) \right) \mathbb{I}(h_l \in A_s) \left[ 1 - \prod_{i=1}^{k-1} \frac{1 - h_{l,i}}{1 + h_{l,i}} \right] \right\rangle_T \\ &= \frac{2 \left( \varphi^{(1)}, \mathbb{T}_{\langle h_l \rangle_T} \varphi^{(1)}(\mathbf{s}) \right)}{\|\varphi\|^2}. \end{aligned}$$

Therefore, for a random tCSP  $(\alpha, p)$  with root degree  $\eta_0 = 1$ , we obtain after averaging, that

$$\widehat{h}_{l+1}^{\text{ave}} \leq 2\mathbb{E}_\varphi \frac{\left(\varphi^{(1)}, \mathbb{T}_{h_l^{\text{ave}}} \varphi^{(1)}\right)}{\|\varphi\|^2},$$

which corresponds to the last inequality of the lemma.  $\square$

We now return to completing the proof of Theorem 4.1.

**Proof (Theorem 4.1, lower bound).** If  $\theta = 1$ ,  $\mathbb{T}_1$  is the identity operator whence  $(\varphi^{(1)}, \mathbb{T}_1 \varphi^{(1)}) = \mathbb{I}_1(\varphi)$ . We have therefore  $F_k(1) = 1/\Omega_k$ . Now, expanding in Fourier series we get,

$$(\varphi^{(1)}, \mathbb{T}_\theta \varphi^{(1)}) = \sum_{Q \subseteq [k-1]} \left| \left(\varphi^{(1)}, \gamma_Q\right) \right|^2 \theta^{|Q|} = \sum_{Q \subseteq [k], Q \ni \{i\}} \left| \left(\varphi, \gamma_Q\right) \right|^2 \theta^{|Q|-1}.$$

By the *Fourier expansion condition*,

$$F_k(\theta) \leq e^{-Ck(1-\theta)}/\Omega_k. \quad (23)$$

Now fix  $\alpha = (1 - \delta)(\Omega_k \log k)/k$ , whence, by Lemma 4.3,  $h_1^{\text{ave}} \leq 1 - k^{-1+\delta}$ , and  $h_1$  is supported on non-negative reals. Using Eq. (19), we get  $\widehat{h}_2^{\text{av}} \leq e^{-Ck^\delta}/\Omega_k$ , and therefore,

$$h_2^{\text{av}} \leq 1 - \exp\{-2(1 - \delta)e^{-Ck^\delta} \log k\} \leq e^{-Ck^\delta/2}.$$

On the other hand, from the Eq. (7), we obtain the following bounds for  $F_k(\theta)$ ,  $R_k(\theta)$ :

$$F_k(\theta) \leq 2\mathbb{E}_\varphi \left[ \frac{\sum_{i=1}^{k-1} \left| \left(\varphi^{(1)}, \gamma_{\{i\}}\right) \right|^2}{\|\varphi\|^2} \right] \theta + 2\mathbb{E}_\varphi \left[ \frac{\mathbb{I}_1(\varphi)}{\|\varphi\|^2} \right] \theta^2 \leq \frac{Ae^{-Ck/2}\theta + \theta^2}{\Omega_k},$$

$$\begin{aligned}
R_k(\theta) &\leq 2\mathbb{E}_{\varphi} \left[ \frac{2\mathbf{I}_1(\varphi)}{\|\varphi\|^2} \sum_{i=1}^{k-1} |(\varphi^{(1)}, \gamma_{\{i\}})|^2 \right] \theta^2 + 2\mathbb{E}_{\varphi} \left[ \frac{2\mathbf{I}_1(\varphi)}{\|\varphi\|^2} \sum_{Q \subseteq [k-1]} |(\varphi^{(1)}, \gamma_Q)| \right] \theta^3 \\
&\leq \frac{Ae^{-Ck/2}\theta^2 + k^a\theta^3}{\Omega_k}.
\end{aligned}$$

Therefore, for all  $\ell$  we have

$$h_{\ell+1}^{\text{av}} \leq 1 - e^{-k\alpha[F_k(h_{\ell}^{\text{av}}) + R_k(h_{\ell}^{\text{av}})]} \leq (1 - \delta) \log k (2Ae^{-Ck/2}h_{\ell}^{\text{av}} + 2k^a(h_{\ell}^{\text{av}})^{3/2}).$$

which implies  $h_{\ell}^{\text{av}} \rightarrow 0$  if, for some  $\ell > 0$ ,  $h_{\ell}^{\text{av}} \leq k^{-5a}$ , thus finishing the proof.  $\square$

## 5 Reconstruction on Trees to Graphs: the case of proper $q$ colorings

In this section we prove that the set of solutions of the proper  $q$ -coloring ensemble satisfies the *sphericity* condition described in Section 3.3.

Given two assignments  $\underline{x}^{(1)}, \underline{x}^{(2)}$  of the variables  $x_1, \dots, x_n$ , their joint type  $v_{\underline{x}^{(1)}, \underline{x}^{(2)}}$  is the  $q \times q$  matrix with  $v_{\underline{x}^{(1)}, \underline{x}^{(2)}}(i, j) \stackrel{\text{def}}{=} \frac{1}{n} \# \{t \in G : \underline{x}^{(1)}(t) = i \text{ and } \underline{x}^{(2)}(t) = j\}$ . We consider random assignments  $\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}$  taken uniformly and independently over all the satisfying assignments of a random instance of the  $q$ -coloring model with edge-variable density  $\alpha$ . Our purpose is to prove that for all  $\delta > 0$ ,  $\|v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v}\|_{\text{TV}} \leq \delta$  w.h.p., where  $\bar{v}$  is the matrix with all entries equal to  $1/q^2$ .

Our argument makes crucial use of the following estimate for the partition function from [AC08].

**Lemma 5.1 ([AC08, Lemma 7])** *Let  $Z$  be the number of satisfying assignments of a random instance of the  $q$ -coloring model with edge-variable density  $\alpha < q \log q$ , then*

$$\mathbf{E}Z \geq \Omega\left(\frac{1}{n^{(q-1)/2}}\right) \left[q\left(1 - \frac{1}{q}\right)^{\alpha}\right]^n,$$

and, for some function  $f(n)$  of order  $o(n)$ , we have  $\text{Prob}(Z < e^{-f(n)}\mathbf{E}[Z]) \rightarrow 0$  as  $n \rightarrow \infty$ .

Let us introduce some notation. If  $w$  is a vector of length  $q$  and  $v$  is a  $q \times q$  matrix  $v$ , let  $\mathcal{H}$  and  $\mathcal{E}$  denote their entropy and their energy respectively, where

$$\begin{aligned}
\mathcal{H}(v) &= -\sum_{i,j} v(i,j) \log v(i,j), \quad \mathcal{H}(w) = -\sum_i w(i) \log w(i) \\
\mathcal{E}(v) &= \log \left( 1 - \sum_i \left( \sum_j v(i,j) \right)^2 - \sum_j \left( \sum_i v(i,j) \right)^2 + \sum_{i,j} v(i,j)^2 \right), \quad \mathcal{E}(w) = \log \left( 1 - \sum_i w(i)^2 \right)
\end{aligned}$$

Let  $\mathcal{B}_q^{\epsilon}$  consists of all the  $q$ -vectors  $w$  with nonnegative entries such that  $\sum_i w(i) = 1$  and  $\|w - \bar{w}\|^2 > \epsilon$ .

Similarly, let  $\mathcal{B}_{q \times q}^{\delta, \epsilon}$  be the set of all the  $q \times q$  matrices with nonnegative entries such that  $\|(v - \bar{v})\mathbf{1}\|^2 \leq \delta$ ,  $\|\mathbf{1}^t(v - \bar{v})\|^2 \leq \delta$  and  $\|v - \bar{v}\|^2 \geq \epsilon$ .

Our goal in this section is to prove the following theorem.

**Theorem 5.2** Let  $\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}$  be random assignments taken uniformly and independently over all the satisfying assignments of a random instance of the  $q$ -coloring model with edge-variable density  $\alpha$ . If  $\alpha < (q-1) \log(q-1)$ , then for any  $\epsilon > 0$ ,

$$\text{Prob} \left( \left\| v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v} \right\|^2 > \epsilon \right) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

We will present several lemmas before returning to the proof of the Theorem. First we introduce estimations concerning an additive functional depending on the energy and entropy of a vector of length  $q$ .

**Lemma 5.3** If  $w \in \mathcal{B}_q^\epsilon$ , then  $\mathcal{H}(w) + \alpha \mathcal{E}(w) \leq [\mathcal{H}(\bar{w}) + \alpha \mathcal{E}(\bar{w})] - \frac{\alpha \epsilon}{2(1-1/q)}$ .

**Proof.** Notice that  $[\mathcal{H}(\bar{w}) + \alpha \mathcal{E}(\bar{w})] - [\mathcal{H}(w) + \alpha \mathcal{E}(w)] = \alpha \log \left( \frac{1 - \frac{\|\bar{w}\|^2}{q}}{1 - \frac{\|w\|^2}{q}} \right)$ . This quantity is bounded below by  $\alpha \log \left( 1 + \frac{\epsilon}{1-1/q} \right)$ , and therefore by  $\frac{\alpha \epsilon}{2(1-1/q)}$ .  $\square$

**Lemma 5.4** Let  $\underline{\mathbf{x}}$  be a random assignment of the variables taken uniformly over all the satisfying assignments of a random instance of the  $q$ -coloring model with edge-variable density  $\alpha < q \log q$ . Then, for any  $\epsilon > 0$ ,

$$\text{Prob} \left( \left\| w_{\underline{\mathbf{x}}} - \bar{w} \right\|^2 > \epsilon \right) \rightarrow 0 \text{ as } n \rightarrow \infty$$

where  $w$  is the vector with  $q$  entries such that  $w_{\underline{\mathbf{x}}}(i) = \frac{1}{n} \# \{v \in G : \underline{\mathbf{x}}_v = i\}$  and  $\bar{w}$  is the vector with all entries equal to  $1/q$ .

**Proof.** Given a property  $P$ , denote by  $Z(P)$ , the number of satisfying assignments for which  $P$  holds. Choose  $\xi$  such that  $\xi < \frac{\alpha \epsilon}{2(1-1/q)}$ . We have that

$$\text{Prob} \left( \left\| w_{\underline{\mathbf{x}}} - \bar{w} \right\|^2 > \epsilon \right) = \mathbf{E} \left[ Z \left( \left\| w_{\underline{\mathbf{x}}} \right\|^2 > \epsilon + 1/q \right) / Z \right],$$

an expression that we can bound by

$$\frac{\mathbf{E} \left[ Z \left( \left\| w_{\underline{\mathbf{x}}} \right\|^2 > \epsilon + 1/q \right) \right]}{e^{-n\xi} \mathbf{E} [Z]} + \text{Prob} \left( Z < e^{-n\xi} \mathbf{E} [Z] \right).$$

Now, according to the Lemma 5.1,  $\text{Prob} \left( Z < e^{-n\xi} \mathbf{E} [Z] \right) \rightarrow 0$ , and therefore it is enough to show that the term

$$\mathbf{E} \left[ Z \left( \left\| w_{\underline{\mathbf{x}}} \right\|^2 > \epsilon + 1/q \right) \right] / e^{-n\xi} \mathbf{E} [Z] \text{ vanishes.}$$

Denote by  $\mathcal{G}_\epsilon$  the set of all vectors  $\ell$ , with nonnegative integer entries, such that  $\sum_{i=1}^q (\ell_i/n) = 1$  and  $\sum_{i=1}^q (\ell_i/n)^2 > \epsilon + 1/q$ , and denote by  $\Omega_w$  the set of assignments  $\underline{\mathbf{x}}$  such that  $w_{\underline{\mathbf{x}}}$  is equal to the vector  $w$ . Now,

$$\begin{aligned} \mathbf{E} \left[ Z \left( \left\| w_{\underline{\mathbf{x}}} \right\|^2 > \epsilon + 1/q \right) \right] &= \sum_{\ell \in \mathcal{G}_\epsilon} \sum_{\underline{\mathbf{x}} \in \Omega_{\ell/n}} \text{Prob} (\underline{\mathbf{x}} \text{ is a satisfying assignment}) \\ &= \sum_{\ell \in \mathcal{G}_\epsilon} \frac{n!}{\prod_{i=1}^q \ell_i!} \left( \left[ \frac{n}{n-1} \right] \left[ 1 - \sum_{i=1}^q (\ell_i/n)^2 \right] \right)^{\alpha n} \\ &\leq \sum_{\ell \in \mathcal{G}_\epsilon} 3q^{2q} \sqrt{n} \exp \left( n [\mathcal{H}(\ell/n) + c_n \mathcal{E}(\ell/n)] \right) \\ &\leq 3q^{2q} \sqrt{n} |\mathcal{G}_\epsilon| \sup_{\ell \in \mathcal{G}_\epsilon} \left\{ \exp \left( n [\mathcal{H}(\ell/n) + c_n \mathcal{E}(\ell/n)] \right) \right\}. \end{aligned} \tag{24}$$

Here  $|\mathcal{G}_\epsilon|$  is the number of elements of  $\mathcal{G}_\epsilon$ , which is bounded by  $n^q$ . Notice also that if  $\ell \in \mathcal{G}_\epsilon$ , then  $\ell/n \in \mathcal{B}_q^\epsilon$ , so that by Lemma 5.3,

$$\begin{aligned} \mathcal{H}(\ell/n) + \alpha\mathcal{E}(\ell/n) &\leq [\mathcal{H}(j_q) + \alpha\mathcal{E}(j_q)] - \frac{\alpha\epsilon}{2(1-1/q)} \\ &= \log q + \alpha \log(1-1/q) - \frac{\alpha\epsilon}{2(1-1/q)}. \end{aligned} \quad (25)$$

On the other hand by the Lemma 5.1, there is some constant  $C$  such that

$$e^{-n\xi} \mathbf{E}[Z] \geq \frac{C}{n^{(q-1)/2}} e^{-n\xi} \left[ q \left(1 - \frac{1}{q}\right)^\alpha \right]^n. \quad (26)$$

Combining Eq. (24), (25) and (26), we have that for a polynomial  $p(n)$  of degree  $3q/2$ ,

$$\frac{\mathbf{E}\left[Z\left(\|w_{\underline{x}} - \bar{w}\|^2 > \epsilon\right)\right]}{e^{-n\xi} \mathbf{E}[Z]} \leq p(n) \exp\left(n\left[\xi - \frac{\alpha\epsilon}{2(1-1/q)}\right]\right). \quad (27)$$

From (27), it is now clear that  $\frac{\mathbf{E}\left[Z\left(\|w_{\underline{x}} - \bar{w}\|^2 > \epsilon\right)\right]}{e^{-n\xi} \mathbf{E}[Z]} \rightarrow 0$  as  $n \rightarrow \infty$ , due to the fact that  $\xi - \frac{\alpha\epsilon}{2(1-1/q)} < 0$ .  $\square$

Next, our objective is to work with the quantity  $\kappa_q^{\delta,\epsilon}$ , which we define as the upper limit of the interval (indeed, easy to see that this is an interval) consisting of the values  $c$  such that

$$\sup_{v \in \mathcal{B}_{q \times q}^{\delta,\epsilon}} \mathcal{H}(v) + c\mathcal{E}(v) \leq \mathcal{H}(\bar{v}) + \alpha\mathcal{E}(\bar{v}).$$

To motivate, let us recall that an important part of the second moment argument of Achlioptas and Naor [AN05, Theorem 7] (in showing that the chromatic number  $\chi[G(n, d/n)]$  concentrated on two possible values), relied on an optimization of the expression  $\mathcal{H}(v) + \alpha\mathcal{E}(v)$  over the Birkoff polytope  $\mathcal{B}_{q \times q}$  of the  $q \times q$  doubly stochastic matrices. In particular, they proved that, as long as  $\alpha \leq (q-1) \log(q-1)$ , one has

$$\sup_{v \in \mathcal{B}_{q \times q}} \mathcal{H}(v) + \alpha\mathcal{E}(v) = \mathcal{H}(\bar{v}) + \alpha\mathcal{E}(\bar{v}). \quad (28)$$

Since  $\mathcal{B}_{q \times q}^{0,\epsilon} \subseteq \mathcal{B}_{q \times q}$ , we have  $\kappa_q^{0,\epsilon} \geq (q-1) \log(q-1)$ . The next lemma says that  $\sup_{v \in \mathcal{B}_{q \times q}^{\delta,\epsilon}} \mathcal{H}(v) + \alpha\mathcal{E}(v)$  is in fact ‘separated’ from  $\mathcal{H}(\bar{v}) + \alpha\mathcal{E}(\bar{v})$ , provided that  $\alpha < \kappa_q^{\delta,\epsilon}$ .

**Lemma 5.5** *Suppose that  $v \in \mathcal{B}_{q \times q}^{\delta,\epsilon}$  where  $\epsilon > 2\delta$ , then, if  $\alpha < \kappa_q^{\delta,\epsilon}$ , we have that*

$$[\mathcal{H}(v) + \alpha\mathcal{E}(v)] \leq [\mathcal{H}(\bar{v}) + \alpha\mathcal{E}(\bar{v})] - \frac{(\kappa_q^{\delta,\epsilon} - \alpha)}{2(1-1/q)^2} [\epsilon - 2\delta].$$

**Proof.** Indeed,

$$\begin{aligned} [\mathcal{H}(\bar{v}) + \alpha\mathcal{E}(\bar{v})] - [\mathcal{H}(v) + \alpha\mathcal{E}(v)] &= \left[ \mathcal{H}(\bar{v}) + \kappa_q^{\delta,\epsilon} \mathcal{E}(\bar{v}) \right] - \left[ \mathcal{H}(v) + \kappa_q^{\delta,\epsilon} \mathcal{E}(v) \right] + (\kappa_q^{\delta,\epsilon} - \alpha) [\mathcal{E}(v) - \mathcal{E}(\bar{v})] \\ &\geq (\kappa_q^{\delta,\epsilon} - \alpha) \left[ \log \left( 1 + \frac{1}{(1-1/q)^2} \left[ \|v - \bar{v}\|^2 - \|(v - \bar{v})\mathbf{1}\|^2 - \|\mathbf{1}^t(v - \bar{v})\|^2 \right] \right) \right] \\ &\geq \frac{(\kappa_q^{\delta,\epsilon} - \alpha)}{2(1-1/q)^2} [\epsilon - 2\delta]. \end{aligned}$$

$\square$

**Lemma 5.6** Given  $\epsilon > 0$  and  $\alpha < \alpha_q = (q-1) \log(q-1)$ , there exists  $\delta > 0$  such that  $\kappa_q^{\delta, \epsilon} \geq \alpha$ .

**Proof.** Assume the contrary, then there exists a sequence  $\delta_n \downarrow 0$  such that  $\kappa_q^{\delta_n, \epsilon} < \alpha$  for each  $n$ . Due to the continuity of  $\exp(\mathcal{H}(v) + \alpha \mathcal{E}(v))$  in the compact set  $\mathcal{B}_{q \times q}^{\delta, \epsilon}$ , the supremum of  $\exp(\mathcal{H}(v) + \alpha_q \mathcal{E}(v))$  is reached at a matrix  $v_{\delta_n} \in \mathcal{B}_{q \times q}^{\delta_n, \epsilon} \subseteq \mathcal{P}_{q \times q}$ , and due to the compactness of  $\mathcal{P}_{q \times q}$ , a subsequence  $\{v_{\delta_{n_k}}\}_{k \geq 1}$  of these matrices converges in  $\mathcal{P}_{q \times q}$  to a matrix  $v \in \mathcal{B}_{q \times q}^{0, \epsilon}$ . Therefore  $\mathcal{H}(v) + \alpha \mathcal{E}(v) \leq \mathcal{H}(\bar{v}) + \alpha \mathcal{E}(\bar{v}) - \frac{(\alpha_q - \alpha)\epsilon}{2(1-1/q)^2}$ . On the other hand,

$$\mathcal{H}(v) + \alpha \mathcal{E}(v) \geq \liminf_{k \rightarrow \infty} \mathcal{H}(v_{\delta_{n_k}}) + \alpha \mathcal{E}(v_{\delta_{n_k}}) \geq \mathcal{H}(\bar{v}) + \alpha \mathcal{E}(\bar{v}),$$

obtaining a contradiction.  $\square$

**Proof of Theorem 5.2.** Given a property  $P$ , denote by  $Z^{(2)}(P)$ , the number of pairs of satisfying assignments for which  $P$  holds. Take  $\alpha'$  such that  $\alpha < \alpha' < (q-1) \log(q-1)$  and use Lemma 5.6 to choose  $\delta$  such that  $\kappa_q^{\delta, \epsilon} \geq \alpha'$ , guaranteeing also that  $2\delta < \epsilon$ . Now, let  $\xi$  be a positive real such that  $2\xi < \frac{(\alpha' - \alpha)}{2(1-1/q)^2} [\epsilon - 2\delta]$ . We have that

$$\text{Prob} \left( \left\| v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v} \right\|^2 > \epsilon \right) = \mathbf{E} \left[ Z^{(2)} \left( \left\| v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v} \right\|^2 > \epsilon \right) / Z^2 \right],$$

which is bounded by the addition of the terms  $E \left[ Z^{(2)} \left( v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} \in \mathcal{B}_{q \times q}^{\delta, \epsilon} \right) \right] / e^{-2n\xi} \mathbf{E} [Z]^2$ ,  $\text{Prob} (Z < e^{-n\xi} \mathbf{E} [Z])$ ,  $\text{Prob} \left( \left\| \left( v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v} \right) \mathbf{1} \right\|^2 > \epsilon \right)$  and  $\text{Prob} \left( \left\| \mathbf{1}^t \left( v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v} \right) \right\|^2 > \epsilon \right)$ . Now, Lemma 5.1 implies that the second term vanishes and lemma 5.4 implies that the last two terms go to zero. Therefore, to show that  $\text{Prob} \left( \left\| v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} - \bar{v} \right\|^2 > \epsilon \right) \rightarrow 0$  is sufficient to prove that the term  $\mathbf{E} \left[ Z^{(2)} \left( v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} \in \mathcal{B}_{q \times q}^{\delta, \epsilon} \right) \right] / e^{-2n\xi} \mathbf{E}$  vanishes.

Denoting by  $\mathcal{G}_{\epsilon, \delta}$  the set of all  $q \times q$  matrices  $L$ , with nonnegative integer entries, such that  $L/n \in \mathcal{B}_{q \times q}^{\delta, \epsilon}$ , and denoting by  $\Omega_v$  the set of pairs of colorings  $x_1, x_2$  such that  $v_{x_1, x_2}$  is equal to the matrix  $v$ , we have

$$\begin{aligned} \mathbf{E} \left[ Z^{(2)} \left( v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} \in \mathcal{B}_{q \times q}^{\delta, \epsilon} \right) \right] &= \sum_{L \in \mathcal{G}_{\epsilon, \delta}} \sum_{x_1, x_2 \in \Omega_{L/n}} \text{Prob} (x_1 \text{ and } x_2 \text{ are satisfying assignments}) \\ &= \sum_{L \in \mathcal{G}_{\epsilon, \delta}} \frac{n!}{\prod_{i,j} L_{ij}!} \left[ \frac{n}{n-1} \right]^{an} \left( 1 - \sum_i \left( \sum_j L_{ij}/n \right)^2 - \sum_j \left( \sum_i L_{ij}/n \right)^2 + \sum_{i,j} (L_{ij}/n)^2 \right)^{an} \\ &\leq \sum_{L \in \mathcal{G}_{\epsilon, \delta}} 3q^{2q} \sqrt{n} \exp (n [\mathcal{H}(L/n) + \alpha E(L/n)]). \end{aligned}$$

And now, because  $\kappa_q^{\delta, \epsilon} \geq \alpha' > \alpha$  and  $L/n \in \mathcal{B}_{q \times q}^{\delta, \epsilon}$  where  $2\delta < \epsilon$ , we can invoke Lemma 5.5 to get that

$$[\mathcal{H}(L/n) + \alpha \mathcal{E}(L/n)] \leq [\mathcal{H}(\bar{v}) + \alpha \mathcal{E}(\bar{v})] - \frac{(\alpha' - \alpha)}{2(1-1/q)^2} [\epsilon - 2\delta].$$

Therefore,

$$\mathbf{E} \left[ Z^{(2)} \left( v_{\underline{\mathbf{x}}^{(1)}, \underline{\mathbf{x}}^{(2)}} \in \mathcal{B}_{q \times q}^{\delta, \epsilon} \right) \right] \leq 3q^{2q} \sqrt{n} |\mathcal{G}_{\epsilon, \delta}| [q(1-1/q)^\alpha]^{2n} \exp \left( -n \frac{(\alpha' - \alpha)}{2(1-1/q)^2} [\epsilon - 2\delta] \right),$$

where  $|\mathcal{G}_{\epsilon,\delta}|$  is the number of elements in  $\mathcal{G}_{\epsilon,\delta}$ , which is bounded by  $n^{q^2}$ . On the other hand by Lemma 5.1, we have that for some constant  $C$ ,

$$e^{-2n\xi} \mathbf{E}[Z]^2 \geq \frac{C}{n^{(q-1)}} e^{-2n\xi} \left[ q \left( 1 - \frac{1}{q} \right)^\alpha \right]^{2n}.$$

Hence, for a polynomial  $p(n)$  of degree  $q^2 + q - 1$ , we have

$$\frac{\mathbf{E} \left[ Z^{(2)} \left( v_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}} \in \mathcal{B}_{q \times q}^{\delta, \epsilon} \right) \right]}{e^{-2n\xi} \mathbf{E}[Z]^2} \leq p(n) \exp \left\{ n \left( 2\xi - \frac{(\alpha' - \alpha)}{2(1 - 1/q)^2} [\epsilon - 2\delta] \right) \right\}.$$

Due to the fact that  $2\xi < \frac{(\alpha' - \alpha)}{2(1 - 1/q)^2} [\epsilon - 2\delta]$ , it is now clear that  $\frac{\mathbf{E} \left[ Z^{(2)} \left( v_{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}} \in \mathcal{B}_{q \times q}^{\delta, \epsilon} \right) \right]}{e^{-2n\xi} \mathbf{E}[Z]^2} \rightarrow 0$  as  $n \rightarrow \infty$ .  $\square$

*Acknowledgments.* The last two authors are grateful to Eric Vigoda and Linji Yang for many insightful discussions on reconstruction problems, and for their role in the early development of this project. The authors also gratefully acknowledge the support and the hospitality of BIRS (Canada) and DIMACS (USA), which provided ideal environs for carrying out a significant part of this research collaboration. Finally, the authors thank an anonymous referee, of a previous version of the paper, for several helpful remarks which resulted in an improved presentation.

## References

- [AC08] D. Achlioptas and A. Coja-Oghlan, *Algorithmic Barriers from Phase Transitions*, Proc. of IEEE FOCS 2008.
- [AM02] D. Achlioptas and C. Moore, *The asymptotic order of the random  $k$ -SAT threshold*, Proc. of IEEE FOCS 2002.
- [AN05] D. Achlioptas and A. Naor, *The two possible values of the chromatic number of a random graph*, Annals of Mathematics, **162** (2005), 1333–1349.
- [ANP05] D. Achlioptas, A. Naor, and Y. Peres, *Rigorous location of phase transitions in hard optimization problems*, Nature **435** (2005), 759–764.
- [AR06] D. Achlioptas and F. Ricci-Tersenghi, *On the solution-space geometry of random constraint satisfaction problems*, Proc. of ACM STOC 2006.
- [AS04] D. Aldous, J. M. Steele, ‘The Objective Method: Probabilistic Combinatorial Optimization and Local Weak Convergence,’ in *Probability on discrete structures*, H. Kesten (ed.), New York, 2004.
- [Bec75] W. Beckner. *Inequalities in Fourier Analysis*. Ann. of Math., **102** (1975), 159–182.
- [BK+05] N. Berger, C. Kenyon, E. Mossel and Y. Peres, *Glauber dynamics on trees and hyperbolic graphs*, Probab. Theory Relat. Fields, **131** (2005) 311-340.
- [BVV07] N. Bhatnagar, J. Vera, and E. Vigoda. *Reconstruction for colorings on trees*. <http://front.math.ucdavis.edu/0711.3664>, 2007.

- [BV04] A. Braunstein and R. Zecchina, “Survey propagation as local equilibrium equations,” *J. of Stat. Mech.* (2004), P06007
- [Bon70] A. Bonami, *Études des coefficients Fourier des fonctions de  $L^p(G)$* . *Ann. Inst. Fourier*, **20** (1970), 335–402.
- [CD+03] S. Cocco, O. Dubois, J. Mandler, R. Monasson, *Rigorous Decimation-Based Construction of Ground States for Spin-Glass Models on Random Lattices*, *Phys. Rev. Lett.* **90** (2003), 047205.
- [CD02] N. Creignou and H. Daude. *Random generalized satisfiability problems*, Proceedings of SAT, Cite-seer, (2002).
- [CD04] N. Creignou and H. Daude. *Combinatorial sharpness criterion and phase transition classification for random CSPs*, *Information and Computation* 190, No. 2 (2004), 220-238.
- [CD09] N. Creignou and H. Daude. *The SAT–UNSAT transition for random constraint satisfaction problems*, *Discrete mathematics*, 309, No 8 (2009), 2085-2099.
- [Dia88] P. Diaconis, *Group representations in probability and statistics. Institute of Mathematical Statistics Lecture Notes*, **11** (1988), Hayward, CA.
- [ES82] P. Erdős and M. Simonovits, *Supersaturated graphs and hypergraphs*, *Combinatorica* 3 (1982), 181–192.
- [F05] E. Friedgut, *Hunting for sharp thresholds*, *Random Structures Algorithms* 26 (2005), 37–51.
- [Geo88] H.-O. Georgii. ‘Gibbs Measures and Phase Transitions,’ de Gruyter, Berlin, 1988.
- [GM07] A. Gerschenfeld, A. Montanari. *Reconstruction for models on random graphs*, Proc. of IEEE FOCS 2007.
- [HPT08] J. Hartigan, D. Pollard, S. Tatikonda. *Conditioned Poisson Distributions and the concentration of chromatic numbers*. <http://www.stat.yale.edu/~pollard/Papers/chromatic.30june08.pdf>
- [KM+07] F. Krzakala, A. Montanari, F. Ricci-Tersenghi, G. Semerjian and L. Zdeborova, *Gibbs States and the Set of Solutions of Random Constraint Satisfaction Problems*, *Proc. Natl. Acad. Sci.* 1004, 10318 (2007)
- [MMW07] E. Maneva, E. Mossel and M. J. Wainwright, *A new look at survey propagation and its generalizations*, *Journal of the ACM* **54** (2007), 17
- [MPZ02] M. Mézard, G. Parisi and R. Zecchina *Analytic and Algorithmic Solution of Random Satisfiability Problems*, *Science* **297** (2002), 812-815.
- [MRZ03] M. Mézard, F. Ricci-Tersenghi and R. Zecchina *Alternative solutions to diluted p-spin models and XORSAT problems*, *J. Stat. Phys.* **111** (2003), 505-53
- [MZ02] M. Mézard and R. Zecchina *Random K-satisfiability problem: From an analytic solution to an efficient algorithm*, *Phys. Rev. E* **66** (2002), 056126
- [MM06] M. Mézard and A. Montanari, *Reconstruction on Trees and Spin Glass Transition*, *J. Stat. Phys.* **124** (2006), 1317-1350.
- [MM09] M. Mézard and A. Montanari, ‘Information, Physics, and Computation,’ Oxford University Press, Oxford, 2009.



- [MMZ05] M. Mézard, T. Mora and R. Zecchina, *Clustering of Solutions in the Random Satisfiability Problem*, Phys. Rev. Lett. **94** (2005), 197205.
- [M05] M. Molloy, *Cores in random hypergraphs and Boolean formulas*, Random Structures and Algorithms **27**, No 1, (2005).
- [MP03] E. Mossel and Y. Peres, *Information Flow on Trees*, Ann. Appl. Probab. **13** (2003), 817-844.
- [Odo08] R. O’Donnell, Some topics in analysis of Boolean functions, FOCS tutorial (2008), and a survey article: found at <http://www.cs.cmu.edu/~odonnell/papers/analysis-survey.pdf>
- [Pa02] G. Parisi, *On local equilibrium equations for clustering states*, arXiv:cs/0212047 (2002)
- [RU08] T. Richardson and R. Urbanke, ‘Modern Coding Theory’, Cambridge University Press, Cambridge, UK, 2008.
- [Sem08] G. Semerjian. *On the freezing of variables in random constraint satisfaction problems*. J. Stat. Phys., 130-251, 2008.
- [Sly08] A. Sly. *Reconstruction of random colourings*. <http://front.math.ucdavis.edu/0802.3487>, 2008.

## A Proof of Proposition 3.1

Given a random instance from the ensemble  $\text{CSP}(n, p, \alpha)$ , let  $\{\varphi_a\}_{a=1}^{\alpha n}$  be its set of clauses and consider the symmetrized statistic

$$L_n(\varphi) = \frac{1}{n\alpha k!} \sum_{\sigma \in S_k} \# \{ a \in [n\alpha] : \varphi_a = \varphi^\sigma \}. \quad (29)$$

It is convenient to introduce two slightly modified ensembles. We denote by  $\text{CSP}(n, p, \alpha; \tilde{p}_n)$  the ensemble  $\text{CSP}(n, p, \alpha)$  conditioned on  $L_n = \tilde{p}_n$ .

A binary configuration  $\underline{x}$  is said to be balanced if  $|\underline{x} \cdot \underline{1}| \leq 1$ . We will use  $Z$  and  $Z_b$ , to denote the variable that counts the number of satisfying assignments and balanced satisfying assignments, respectively, of a random CSP ensemble. Given two binary assignments  $\underline{x}^{(1)}, \underline{x}^{(2)}$ , we define their overlap as

$$Q_{12} \stackrel{\text{def}}{=} \frac{1}{n} \underline{x}^{(1)} \cdot \underline{x}^{(2)} = \frac{1}{n} \sum_{i=1}^n x_i^{(1)} x_i^{(2)}. \quad (30)$$

In other words  $(1 - Q_{12})/2$  is the normalized Hamming distance of  $\underline{x}^{(1)}$  and  $\underline{x}^{(2)}$ .

**Proof (Proposition 3.1, upper bound).** The upper bound in Proposition 3.1 follows from a first moment calculation. Let  $Z$  be the number of solutions of a random instance from the ensemble  $\overline{\text{CSP}}(n, p, \alpha)$ . We will show that, for  $\alpha > (1 + \epsilon)\hat{\Omega}_k \log 2$ ,  $\mathbf{E}[Z] \rightarrow 0$  as  $n \rightarrow \infty$ . First fix  $\tilde{p}_n$  such that  $\|\tilde{p}_n - p\|_{\text{TV}} \leq 1/n^{1/2-\gamma}$ . Notice that the probability that a random clause of type  $\varphi$  is satisfied by the assignment  $x$  with

$x \cdot 1 = n\theta$  is  $\|\varphi\|_\theta^2$ . This implies

$$\begin{aligned}
\mathbf{E}[Z|L_n = \tilde{p}_n] &= \sum_{x \in \{-1,1\}^n} \mathbf{P}(x \text{ is a satisfying assignment} | L_n = \tilde{p}_n) \\
&\leq n \sup_{\theta \in [-1,1]} \sum_{x \cdot 1 = n\theta} \mathbf{P}(x \text{ is a satisfying assignment} | L_n = \tilde{p}_n) \\
&\leq n 2^n \prod_{\varphi} \|\varphi\|_\theta^{2\tilde{p}_n(\varphi)\alpha n} \\
&\leq n \exp \left( n \left\{ \log 2 + \alpha \sum_{\varphi} p(\varphi) \log \|\varphi\|_\theta^2 + O(n^{-1/2+\gamma}) \right\} \right) \\
&\leq n \exp \left( n \left\{ \log 2 + \alpha \sum_{\varphi} p(\varphi) \log \|\varphi\|^2 + O(n^{-1/2+\gamma}) \right\} \right),
\end{aligned}$$

where in the last step we used the condition of dominance of balanced assignments. By taking expectation over  $\tilde{p}_n$ , we obtain  $\mathbf{E}[Z] \rightarrow 0$  whenever  $\alpha > (1 + \epsilon) \hat{\Omega}_k \log 2$ , as claimed.  $\square$

To establish the corresponding lower bound, we use the second moment method, but first we need a few preliminary lemmas.

We define by  $\mathcal{K}_n(p; a, A, \gamma)$  to be the set of probability distributions  $\{\tilde{p}(\varphi)\}$  over clauses  $\varphi : \{+1, -1\} \rightarrow \{0, 1\}$  such that:

- (i)  $\text{supp}(\tilde{p}) = \text{supp}(p)$ ;
- (ii)  $\tilde{p}$  satisfies conditions 1-4 and **(a)**, **(b)** stated in Section 2, with constants  $a, A$ ; and finally,
- (iii)  $\|\tilde{p} - p\|_{\text{TV}} \leq n^{-1/2+\gamma}$  for some  $\gamma > 0$ .

Then we have the following.

**Lemma A.1** *Let  $L_n$  be the statistics defined in Eq. (29) for a random formula from the CSP( $n, p, \alpha$ ) ensemble. Then there exists constants  $a, A$  such that for any  $\gamma > 0$ , with high probability*

$$L_n \in \mathcal{K}_n(p; a, A, \gamma). \quad (31)$$

**Proof.** Notice that for each permutation  $\pi$   $L_n(\varphi^\pi) = L_n(\varphi)$  and that, for each  $\varphi \in \{-1, +1\}^k \rightarrow \{0, 1\}$ ,  $k!L_n(\varphi)$  is distributed as a binomial with parameters  $n\alpha$ , and  $k!p(\varphi)$ . In particular  $L_n(\varphi) = 0$  if  $p(\varphi) = 0$  and  $L_n(\varphi) > 0$  with high probability otherwise. This implies Item (i) in the definition of  $\mathcal{K}_n(p; a, A, \gamma)$ .

Item (iii) that,  $\|L_n - p\|_{\text{TV}} \leq n^{-1/2+\gamma}$ , follows immediately from the central limit theorem.

Consider finally Item (ii). Condition 1 is enforced by the symmetrization procedure in Eq. (29). Conditions 2, 3 only depend on  $\text{supp}(L_n)$  and thus hold with high probability by the above argument.

Dominance of balanced assignments (condition 4) is the statement that

$$\mathbb{E}_\varphi \log \|\varphi\|_\theta - \mathbb{E}_\varphi \log \|\varphi\| < 0, \quad (32)$$

for all  $\theta \neq 0$ ,  $\theta \in [-1, 1]$ . Notice that the left hand side is a polynomial in  $\theta$  whose coefficients are continuous function of the quantities  $\{L_n(\varphi)\}$ . Hence this condition is of the form  $L_n \in \mathcal{A}$  for  $\mathcal{A}$  an open set in  $\mathbb{R}^D$ ,  $D = 2^{2^k}$ . Since  $p \in \mathcal{A}$  and  $\|L_n - p\|_{\text{TV}} \leq n^{-1/2+\gamma}$  whp, we conclude  $L_n \in \mathcal{A}$ .

Finally conditions **(a)** and **(b)** only depend on  $\text{supp}(L_n)$  and therefore follow from the above.  $\square$

**Lemma A.2** Given  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ , consider a random instance from the  $\text{CSP}(n, p, \alpha; \tilde{p}_n)$  ensemble. For  $\theta \in \{-1, -1 + 2/n, \dots, 1 - 2/n, 1\}$ , let  $Z_b(Q_{12} = \theta)$  be the number of balanced solution pairs  $\underline{x}^{(1)}, \underline{x}^{(2)} \in \{+1, -1\}^n$  with overlap  $\theta$ . Then,

$$\frac{\mathbf{E} [Z_b(Q_{12} = \theta)]}{[\mathbf{E} Z_b]^2} \leq C n^{-1/2} \exp \{n \Phi(\theta)\},$$

where  $C$  is bounded uniformly in  $\theta$  and

$$\Phi(\theta) \stackrel{\text{def}}{=} H(\theta) + \alpha \mathbb{E}_{\varphi \sim \tilde{p}_n} \log \left\{ \frac{(\varphi, T_\theta \varphi)}{\|\varphi\|^4} \right\}.$$

Here  $H(\theta) \equiv -\frac{1+\theta}{2} \log(1+\theta) - \frac{1-\theta}{2} \log(1-\theta)$  is the binary entropy function.

**Proof.** For simplicity take  $n$  to be even (the argument is analogous for  $n$  odd). Let  $\varphi$  be a boolean function, and let  $i : [k] \rightarrow [n]$  be a uniform random choice of the indexes of the variables in  $\varphi$  (i.e.  $i(1), \dots, i(k)$  are independent and uniform in  $[n]$ ). Given two *balanced* vectors  $\underline{x}^{(1)}, \underline{x}^{(2)} \in \{+1, -1\}^n$ , with  $Q_{12} = \theta$ , we have

$$\mathbf{E}_\pi \left[ \varphi \left( x_{i(1)}^{(1)}, \dots, x_{i(k)}^{(1)} \right) \varphi \left( x_{i(1)}^{(2)}, \dots, x_{i(k)}^{(2)} \right) \right] = (\varphi, T_\theta \varphi).$$

Therefore

$$\begin{aligned} \mathbf{E} Z_b(|Q_{12}| = \theta) &= \sum_{\underline{x}^{(1)}, \underline{x}^{(2)} = n\theta} \mathbf{P} \left( \underline{x}^{(1)}, \underline{x}^{(2)} \text{ are satisfying assignments} \right) \\ &\leq \sum_{\underline{x}^{(1)}, \underline{x}^{(2)} = n\theta} \prod_{\varphi} (\varphi, T_\theta \varphi)^{\tilde{p}_n(\varphi)n\alpha} \\ &\leq \frac{C}{n^{3/2}} \exp \left( n \left\{ \mathcal{H} \left( \frac{1+\theta}{4}, \frac{1+\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4} \right) + \alpha \sum_{\varphi} \tilde{p}_n(\varphi) \log (\varphi, T_\theta \varphi) \right\} \right). \end{aligned}$$

Where  $\mathcal{H}$  is the entropy function

$$\mathcal{H}(\theta_1, \dots, \theta_d) = - \sum_{i=1}^d \theta_i \log \theta_i. \quad (33)$$

and we used the following bound on binomial coefficients (valid for  $\theta_i \geq 0, \theta_1 + \dots + \theta_d = 1$ )

$$\frac{n!}{\prod_{i=1}^d (n\theta_i!)} \leq \frac{C}{n^{(d-1)/2}} \exp \{ \mathcal{H}(\theta_1, \dots, \theta_d) \}. \quad (34)$$

By the very same argument, for some positive  $C'$ ,

$$\begin{aligned} \mathbf{E} Z_b &= \sum_{\underline{x} \text{ balanced}} \prod_{\varphi} \|\varphi\|^{2\tilde{p}_n(\varphi)\alpha n} \\ &> \frac{C'}{n^{1/2}} \exp \left( n \left\{ \mathcal{H} \left( \frac{1}{2}, \frac{1}{2} \right) + \alpha \sum_{\varphi} \tilde{p}_n(\varphi) \log \|\varphi\|^2 \right\} \right). \end{aligned}$$

It is straightforward now to check that

$$\frac{\mathbf{E} Z_b(Q_{12} = \theta)}{(\mathbf{E} Z_b)^2} \leq \frac{C''}{n^{1/2}} \exp \{n \Phi(\theta)\} \quad (35)$$

which implies the claim.  $\square$

**Lemma A.3** Given  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ , consider a random instance from the  $\text{CSP}(n, p, \alpha; \tilde{p}_n)$  ensemble, and define

$$\Omega_{k, \tilde{p}_n} \stackrel{\text{def}}{=} \mathbb{E}_{\varphi \sim \tilde{p}_n} \frac{2\mathbf{I}_1(\varphi)}{1 - 2\mathbf{I}_1(\varphi)}. \quad (36)$$

If  $\alpha \leq (1 - \varepsilon)\Omega_{k, \tilde{p}_n} \log 2$ , then exist a constants  $C_0 = C_0(p; a, A, \gamma, \varepsilon) > 0$  (independent of  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ ), and an absolute constant  $C$  such that for any  $\theta \in \{-1, -1 + 2/n, \dots, 1 - 2/n, 1\}$

$$\frac{\mathbf{E}[Z_b(Q_{12} = \theta)]}{(\mathbf{E}Z_b)^2} \leq \frac{C}{n^{1/2}} e^{-nC_0\theta^2}. \quad (37)$$

**Proof.** In view of the previous lemma, it is sufficient to prove that there exist a constant  $C_0 = C_0(p; a, A, \gamma, \varepsilon) > 0$  (independent of  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ ) such that

$$\Phi(\theta) \leq -C_0\theta^2. \quad (38)$$

Since throughout this proof  $\tilde{p}_n$  is fixed, it will be understood that  $\varphi \sim \tilde{p}_n$  whenever we take expectation over the clause distribution. Also, dependence of  $\Omega_{k, \tilde{p}_n}$  and  $\widehat{\Omega}_{k, \tilde{p}_n}$  (defined analogously) upon  $\tilde{p}_n$  will be dropped.

Fix  $\alpha \leq (1 - \varepsilon)\Omega_k \log 2 \leq (1 - \varepsilon)\widehat{\Omega}_k \log 2$ . We will prove the thesis claim by considering three different regimes for  $\theta$ :  $0 < \theta \leq e^{-ck}$ ,  $e^{-ck} \leq \theta \leq 1 - \varepsilon^{1/2}$  and  $1 - \varepsilon^{1/2} \leq \theta \leq 1$ , where  $c$  is a small constant. In the first two intervals we will prove that the derivative of  $\Phi(\theta)$  with respect to  $\theta$  is strictly negative. Recalling that  $\|\varphi\|^2 \geq 1/2$ , we have

$$\begin{aligned} \frac{d\Phi}{d\theta} &\leq -\text{atanh } \theta + k\alpha \mathbb{E}_\varphi \frac{(\varphi^{(1)}, \mathbf{T}_\theta \varphi^{(1)})}{\|\varphi\|^4} \\ &\leq -\theta + 2k\alpha \mathbb{E}_\varphi \frac{\sum_{i=1}^{k-1} |\varphi_{\{i\}}^{(1)}|^2}{\|\varphi\|^2} \theta + 2k\alpha \mathbb{E}_\varphi \frac{\|\varphi^{(1)}\|^2}{\|\varphi\|^2} \theta^3 \\ &\leq -\theta + Ae^{-Ck} \frac{\alpha}{\Omega_k} \theta + 2k \frac{\alpha}{\Omega_k} \theta^2 \\ &\leq -\frac{1}{2}\theta + 4k\theta^2, \end{aligned}$$

where we used (from Eq. (3)) the hypothesis on low weight Fourier coefficients. The last expression is strictly negative if  $0 < \theta < e^{-ck}$  for any  $c > 0$  and all  $k$  large enough. Integrating the last expression over  $\theta$ , we get  $\Phi(\theta) \leq -C_0\theta^2$ .

Next assume  $e^{-ck} \leq \theta \leq 1 - \varepsilon$ . Using the hypothesis  $(\varphi^{(1)}, \mathbf{T}_\theta \varphi^{(1)}) \leq e^{-Ck(1-\theta)}\|\varphi^{(1)}\|^2$ , we have

$$\begin{aligned} \frac{d\Phi}{d\theta} &\leq -\text{atanh } \theta + 4k\alpha \mathbb{E}_\varphi \frac{\|\varphi^{(1)}\|^2}{\|\varphi\|^4} e^{-Ck\varepsilon} \\ &\leq -\text{atanh } \theta + 2k \frac{\alpha}{\Omega_k} e^{-Ck\sqrt{\varepsilon}} \leq -\text{atanh } \theta + 2(\log 2) k e^{-Ck\varepsilon}, \end{aligned}$$

which is strictly negative if  $\theta > c^{-ak}$  with, say,  $c = (C\varepsilon^2)/2$ .

Finally, we notice that, for  $1 - \varepsilon^2 \leq \theta \leq 1$ , any  $\varepsilon$  small enough we have  $H(\theta) \leq -\log 2 + \varepsilon/10$ . Further, using the fact that  $(\varphi, \mathbf{T}_\theta \varphi) = \|\mathbf{T}_{\theta^{1/2}} \varphi\|^2$  is non-decreasing in  $\theta$

$$\Phi(\theta) \leq -\log 2 + \frac{\varepsilon}{10} - \alpha \mathbb{E}_\varphi \log \|\varphi\|^2 = -\log 2 + \frac{\varepsilon}{10} + \frac{\alpha}{\widehat{\Omega}_k} \leq -\varepsilon \frac{\log 2}{2},$$

which finishes the proof.  $\square$

**Proof (Proposition 3.1, lower bound).** Fix  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ ,  $\alpha \leq (1 - \varepsilon)\Omega_{k, \tilde{p}_n} \log 2$  and let  $Z_b$  the number of balanced solutions of a random instance from the  $\text{CSP}(n, p, \alpha; \tilde{p}_n)$  ensemble. From Lemma A.3 we have that, for  $U_n \equiv \{-1, -1 + 2/n, \dots, 1 - 2/n, 1\}$ ,

$$\frac{\mathbf{E}\{Z_b^2\}}{\{\mathbf{E} Z_b\}^2} = \sum_{\theta \in U_n} \frac{\mathbf{E}\{Z_b(Q_{12} = \theta)\}}{\{\mathbf{E} Z_b\}^2} \quad (39)$$

$$\leq \frac{C}{n^{1/2}} \sum_{\theta \in U_n} e^{-C_0 n \theta^2} \quad (40)$$

$$\leq \frac{C'}{n^{1/2}} n \int_{-\infty}^{\infty} e^{-C_0 n \theta^2} d\theta \leq C'_0 \quad (41)$$

for some new constant  $C'_0 = C'_0(p; a, A, \gamma, \varepsilon) > 0$ .

For  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ , we gave  $\Omega_{k, \tilde{p}_n} = \Omega_k(1 + O(n^{-1/2+\gamma}))$ . Let  $\mathcal{F}_n$  be a random instance from  $\text{CSP}(n, p, \alpha)$  ensemble,  $\tilde{p}_n \in \mathcal{K}_n(p; a, A, \gamma)$ ,  $\alpha \leq (1 - 2\varepsilon)\Omega_k \log 2$ , whence  $\alpha \leq (1 - \varepsilon)\Omega_{k, \tilde{p}_n} \log 2$ . By Paley-Zygmund inequality

$$\mathbf{P}(\mathcal{F}_n \text{ is sat} | L_n = \tilde{p}_n) \geq \frac{\mathbf{E}\{Z_b^2\}}{2\{\mathbf{E} Z_b\}^2} > C'_0/2. \quad (42)$$

By Lemma A.1 we have  $\mathbf{P}(\mathcal{F}_n \text{ is sat}) \geq C'_0/4$ . Finally, the fact that the satisfiability property (of our CSP ensembles) exhibits a sharp transition, thanks to the theorem of N. Creignou and H. Daude [CD09] (see Theorem C.1 in Appendix C here) implies  $\mathbf{P}(\mathcal{F}_n \text{ is sat}) \rightarrow 1$  as  $n \rightarrow \infty$ .  $\square$

## B Proof of Theorem 3.3

In this appendix we introduce the planted CSP ensemble, clarify its connection to the original ensemble, and use it to prove Theorem 3.3. Throughout the section, we denote a CSP instance with  $n\alpha$  clauses by  $F = (F_1, F_2, \dots, F_{n\alpha})$ . Here

$$F_a = (\varphi_a; i_a(1), \dots, i_a(k)) \quad (43)$$

denotes the clause labeled  $a$ , which is completely specified by the Boolean function  $\varphi_a : \{+1, -1\}^k \rightarrow \{0, 1\}$  and by the choice of  $k$  indices  $i_a(1), \dots, i_a(k)$ . The number of solutions of the instance  $F$  is denoted by  $Z(F)$ .

Given a distribution  $p = \{p(\varphi)\}$ , it is also convenient to define the ‘average clause’  $\bar{\varphi} : \{+1, -1\}^n \rightarrow \mathbb{R}_+$ :

$$\bar{\varphi}(\underline{x}) = \frac{1}{n^k} \sum_{i(1), \dots, i(k) \in [n]} \sum_{\varphi} p(\varphi) \varphi(x_{i(1)}, \dots, x_{i(k)}). \quad (44)$$

Throughout this section, we will assume that the strong balance condition (9) holds. We think that this condition can be refined at the price of a more careful analysis.

## B.1 The planted ensemble and a transfer theorem

Given  $n \in \mathbb{N}$ ,  $\alpha \geq 0$ , and a distribution  $p = \{p(\varphi)\}$  over  $k$ -clauses, the planted ensemble  $\text{pCSP}(n, \alpha, p)$  is a joint distribution over binary assignments  $\underline{x}^* = (x_1^*, x_2^*, \dots, x_n^*) \in \{0, 1\}^n$  and random CSP formulas  $F$  defined as follows. The assignment  $\underline{x}^*$  is drawn with distribution

$$\mathbf{P}_p(\underline{x}) \equiv \frac{1}{\mathbf{E} Z(F)} \bar{\varphi}(\underline{x})^{n\alpha}. \quad (45)$$

It is easy to check that this is normalized, i.e. that  $\mathbf{E} Z(F) = \sum_{\underline{x}} \bar{\varphi}(\underline{x})^{n\alpha}$ .

We will use  $\mathbf{P}_p$ ,  $\mathbf{E}_p$  to denote probability and expectation with respect to the planted model. Sampling  $\underline{x}$  from this distribution is straightforward, since  $\mathbf{P}_p(\underline{x})$  is uniform once we condition on the weight of  $\underline{x}$  (i.e. on  $\underline{x} \cdot \underline{1}$ ).

Conditional on  $\underline{x}^*$ , the clauses  $F_a$ ,  $a = 1, 2, \dots, n\alpha$ , are independent and distributed according to

$$\mathbf{P}_p\{F_a = (\varphi_a, i_a(1), \dots, i_a(k)) | \underline{x}^*\} \equiv \frac{1}{n^k \bar{\varphi}(\underline{x}^*)} p(\varphi_a) \varphi_a(x_{i_a(1)}^*, \dots, x_{i_a(k)}^*). \quad (46)$$

draw indices  $i_a(1), \dots, i_a(k) \in [n]$  independently and uniformly at random. Notice that this is indeed a well defined distribution over clauses, and in particular it is normalized thanks to Eq. (44). In order to sample from the above clause distribution, one can proceed as follows. Sample indices  $i_a(1), \dots, i_a(k) \in [n]$  independently and uniformly at random and a boolean function  $\varphi_a$  with distribution  $p(\cdot)$ . If  $\varphi_a(x_{i_a(1)}^*, \dots, x_{i_a(k)}^*) = 1$ , accept this choice, otherwise reject it and repeat the sampling.

The joint distributon of the planted assignment and the CSP instance is then

$$\mathbf{P}_p(F, \underline{x}^*) = \frac{1}{n^{nk\alpha} \mathbf{E} Z(F)} \prod_{a=1}^{n\alpha} p(\varphi_a) \varphi_a(x_{i_a(1)}^*, \dots, x_{i_a(k)}^*). \quad (47)$$

By construction, the assignment  $\underline{x}^*$  satisfies  $F$ . It is convenient to compare the planted distribution with the uniform distribtion we have been considering so far. In this case, an instance is drawn according to the ensemble  $\text{CSP}(n, \alpha, p)$ , and an assignment  $\underline{x}^*$  is drawn uniformly at random from among the ones satisfying  $F$ . The joint distribution is then

$$\mathbf{P}_p(F, \underline{x}^*) = \frac{1}{n^{nk\alpha} Z(F)} \prod_{a=1}^{n\alpha} p(\varphi_a) \varphi_a(x_{i_a(1)}^*, \dots, x_{i_a(k)}^*). \quad (48)$$

By taking the ratio of the above probabilities, we immediately get the following

**Lemma B.1** *Let  $\mathcal{F} : (F, \underline{x}^*) \rightarrow \mathbb{R}$  be a function of an instance, solution pair. Its expectations with respect to the planted and uniform model are related as follows*

$$\mathbf{E}_p \mathcal{F}(F, \underline{x}^*) = \mathbf{E} \left\{ \frac{Z(F)}{\mathbf{E} Z(F)} \mathcal{F}(F, \underline{x}^*) \right\}. \quad (49)$$

**Proof.** By a standard change-of-measure argument  $\mathbf{E}_p \mathcal{F}(F, \underline{x}^*)$  is equal to

$$\begin{aligned} \sum_{(F, \underline{x}^*)} \mathbf{P}_p(F, \underline{x}^*) \mathcal{F}(F, \underline{x}^*) &= \sum_{(F, \underline{x}^*)} \mathbf{P}(F, \underline{x}^*) \left\{ \frac{\mathbf{P}_p(F, \underline{x}^*)}{\mathbf{P}(F, \underline{x}^*)} \mathcal{F}(F, \underline{x}^*) \right\} \\ &= \sum_{(F, \underline{x}^*)} \mathbf{P}(F, \underline{x}^*) \left\{ \frac{Z(F)}{\mathbf{E} Z(F)} \mathcal{F}(F, \underline{x}^*) \right\}, \end{aligned} \quad (50)$$

which is nothing but our claim.  $\square$

It is clear that the planted and uniform model are strictly related as soon as  $Z(F)$  concentrates around its expectation  $\mathbf{E} Z(F)$ .

**Lemma B.2** *Fix  $\alpha < \Omega_k \log 2\{1 + o_k(1)\}$  and let  $Z(F)$  be the number of solutions of a random instance  $F$  from the  $\text{CSP}(n, \alpha, p)$  ensemble. Then, for any  $\varepsilon > 0$ ,  $Z(F) > e^{-n\varepsilon} \mathbf{E} Z(F)$  with high probability.*

**Proof.** For any constant  $A$ , the property  $Z(F) > e^{nA}$  is monotone over the space of CSP instances (regarded as a product space). Applying as in [AC08], a sharp threshold result (which we prove as Lemma C.2 in Appendix C), it is sufficient to prove that  $Z(F) > e^{-n\varepsilon} \mathbf{E} Z(F)$  with probability bounded away from 0 as  $n \rightarrow \infty$ .

Let  $Z_b(F)$  the number of balanced solutions (i.e. the number of solutions such that  $|\underline{x} \cdot \underline{1}| \leq 1$ ). Obviously  $Z(F) \geq Z_b(F)$ . On the other hand, by an argument already employed in Appendix A (here  $U_n \equiv \{-1, -1 + 2/n, \dots, 1 - 2/n, 1\}$ ):

$$\begin{aligned} \mathbf{E}\{Z(F)\} &= \sum_{x \in \{-1, 1\}^k} \mathbf{P}(x \text{ is a satisfying assignment}) \\ &\leq \sum_{\theta \in U_n} \binom{n}{n(1+\theta)/2} \mathbb{E}_\varphi\{\|\varphi\|_\theta\}^2 \\ &\leq \sum_{\theta \in U_n} \binom{n}{n(1+\theta)/2} \mathbb{E}_\varphi\{\|\varphi\|\}^2 \\ &\leq n \binom{n}{n/2} \mathbb{E}_\varphi\{\|\varphi\|\}^2 = n \mathbf{E}\{Z_b(F)\}. \end{aligned}$$

That is  $\mathbf{E}\{Z(F)\}$  and  $\mathbf{E}\{Z_b(F)\}$  differ at most by a polynomial factor. It is therefore sufficient to prove that  $Z_b(F) > e^{-n\varepsilon} \mathbf{E} Z_b(F)$  with probability bounded away from 0 as  $n \rightarrow \infty$ .

This follows from Paley-Zygmund inequality, since

$$\mathbf{P}\left\{Z_b(F) \geq \frac{1}{2} \mathbf{E} Z_b(F)\right\} \geq \frac{\mathbf{E}\{Z_b(F)\}^2}{4\mathbf{E}\{Z_b(F)^2\}} \geq \frac{1}{4C}, \quad (51)$$

for some uniformly bounded  $C > 0$ , by Eq. (41).  $\square$

**Theorem B.3** *Given a sequence of events  $\{A_n\}$  and a constant  $c > 0$ , assume that  $(\underline{x}^*, F) \in A_n$  with probability larger than  $1 - e^{-cn}$  under the planted model  $\text{pCSP}(n, \alpha, p)$ . Then  $(\underline{x}^*, F) \in A$  with high probability under the uniform model.*

**Proof.** Consider the complement of  $A_n$ , denoted by  $A_n^c$ . By Lemma B.1, we have

$$\begin{aligned} \mathbf{P}_p\{(\underline{x}^*, F) \in A_n^c\} &= \mathbf{E}\left\{\frac{Z(F)}{\mathbf{E}Z(F)} \mathbb{I}_{(\underline{x}^*, F) \in A_n^c}\right\} \\ &\geq \mathbf{E}\left\{\frac{Z(F)}{\mathbf{E}Z(F)} \mathbb{I}_{(\underline{x}^*, F) \in A_n^c} \mathbb{I}_{Z(F) \geq e^{-cn/2} \mathbf{E}Z(F)}\right\} \\ &\geq e^{-cn/2} \left\{ \mathbf{P}\{(\underline{x}^*, F) \in A_n^c\} - \mathbf{P}\{(\underline{x}^*, F) \in A_n^c, Z < e^{-cn/2} \mathbf{E}Z(F)\} \right\}. \end{aligned}$$

By solving for  $\mathbf{P}\{(\underline{x}^*, F) \in A_n^c\}$  we get

$$\mathbf{P}\{(\underline{x}^*, F) \in A_n^c\} \leq e^{cn/2} \mathbf{P}_p\{(\underline{x}^*, F) \in A_n^c\} + \mathbf{P}\{Z < e^{-cn/2} \mathbf{E}Z(F)\}.$$

The first term vanishes by assumption, and the second by Lemma B.2.  $\square$

## B.2 Clustering

The proof of Theorem 3.3 proceed in two steps. First we consider a pair  $(\underline{x}^*, F)$  drawn according to the planted model and show that the planted solution is isolated from most of the other solutions. Next, we use Theorem B.3 to transfer this statement to the uniform ensemble.

In order to establish the first result, we need the following estimate.

**Lemma B.4** *Let  $(\underline{x}^*, F)$  be an solution/instance pair distributed according to the planted model, and denote by  $Z^{(2)}(\theta)$  the number of solutions  $\underline{x}$  of  $F$  such that  $\underline{x}^* \cdot \underline{x} = n\theta$ . Then, for any  $a < 1$*

$$\mathbf{E}_p\{Z^{(2)}(\theta) \mid \underline{x}^* \cdot \underline{1} \leq n^a\} = \exp\{n\Psi(\theta) + o(n)\}, \quad (52)$$

$$\Psi(\theta) \equiv H(\theta) + \alpha \log \left\{ \frac{\mathbb{E}_\varphi(\varphi, T_\theta \varphi)}{\mathbb{E}_\varphi\{\|\varphi\|^2\}} \right\}. \quad (53)$$

**Proof.** For the sake of simplicity we shall focus on the case  $\underline{x}^* \cdot \underline{1} = 0$  (i.e.  $n$  is even and the planted solution is perfectly balanced). It should be clear from the derivation that allowing for  $|\underline{x}^* \cdot \underline{1}| \leq n^a$  only produces a change of order  $O(n^{-1+a})$  in the exponent.

Fix such a planted solution  $\underline{x}^*$ , and let  $\underline{x}$  be such that

$$\sum_{i:x_i^*=+1} x_i^* x_i = \frac{n}{2} \theta_+, \quad \sum_{i:x_i^*=-1} x_i^* x_i = \frac{n}{2} \theta_-, \quad (54)$$

with  $(\theta_+ + \theta_-)/2 = \theta$  (whence  $\underline{x}^* \cdot \underline{x} = n\theta$ ). Then

$$\mathbf{P}_p(\underline{x} \text{ is a solution} \mid \underline{x}^*) = [\mathbf{P}_p(\varphi_a(x_{i_a(1)}, \dots, x_{i_a(k)}) = 1 \mid \underline{x}^*)]^{n\alpha}, \quad (55)$$

and by the definition of planted ensemble

$$\begin{aligned} \mathbf{P}_p(\varphi_a(x_{i_a(1)}, \dots, x_{i_a(k)}) = 1 \mid \underline{x}^*) &= \frac{1}{n^k \overline{\varphi}(\underline{x}^*)} \sum_{i_a(1), \dots, i_a(k)} \sum_{\varphi} p(\varphi) \varphi(x_{i_a(1)}^*, \dots, x_{i_a(k)}^*) \varphi(x_{i_a(1)}, \dots, x_{i_a(k)}) \\ &= \frac{1}{\overline{\varphi}(\underline{x}^*)} \mathbb{E}_\varphi(\varphi, S_{\theta_+, \theta_-} \varphi), \end{aligned}$$

where we introduced the operator  $S_{\theta_+, \theta_-}$  acting as follows

$$S_{\theta_+, \theta_-} \varphi(x_1, \dots, x_k) \equiv \sum_{y \in \{+1, -1\}^k} \prod_{i=1}^k \frac{1 + \theta_{x_i} x_i y_i}{2} \varphi(y_1, \dots, y_k). \quad (56)$$

Further

$$\mathbf{P}_p(\underline{x}^* \cdot \underline{1} = 0) = \frac{1}{\mathbf{E}Z(F)} \overline{\varphi}(\underline{x}^*)^{n\alpha} \binom{n}{n/2}.$$

Combining the above, and after a few algebraic manipulations, we get

$$\mathbf{E}_p\{Z^{(2)}(\theta) \mid \underline{x}^* \cdot \underline{1} = 0\} = \frac{1}{\overline{\varphi}(\underline{x}^*)^{n\alpha}} \sum_{\theta_+ + \theta_- = 2\theta} \binom{n/2}{n(1+\theta_+)/4} \binom{n/2}{n(1+\theta_-)/4} [\mathbb{E}_\varphi(\varphi, S_{\theta_+, \theta_-} \varphi)]^{n\alpha},$$



where the sum runs over  $\theta_+, \theta_- \in \{-1, -1 + 4/n, \dots, 1 - 4/n, 1\}$ . Now letting  $\delta = (\theta_+ - \theta_-)/2$  and passing to the Fourier transform, we get

$$\mathbb{E}_\varphi(\varphi, S_{\theta_+, \theta_-} \varphi) = \sum_{Q_1 \subseteq Q_2} \mathbb{E}_\varphi\{\varphi_{Q_1} \varphi_{Q_2}\} \theta^{|Q_1|} \delta^{|Q_2| - |Q_1|} \leq \sum_Q \mathbb{E}_\varphi\{\varphi_Q^2\} \theta^{|Q|} = \mathbb{E}_\varphi(\varphi, S_{\theta, \theta} \varphi),$$

where we used Eq. (10). Also notice that  $(\varphi, S_{\theta, \theta} \varphi) = (\varphi, T_\theta \varphi)$ . Therefore, the sum over  $\theta_+, \theta_-$  can be estimated by the  $\theta_+ = \theta_-$  term, up to a polynomial factor

$$\mathbb{E}_P\{Z^{(2)}(\theta) | \underline{x}^* \cdot \mathbf{1} = 0\} = \frac{1}{\bar{\varphi}(\underline{x}^*)^{n\alpha}} n^{O(1)} \left( \frac{n/2}{n(1+\theta)/4} \right)^2 [\mathbb{E}_\varphi(\varphi, T_\theta \varphi)]^{n\alpha}.$$

The statement follows by noticing that  $\bar{\varphi}(\underline{x}^*) = \mathbb{E}_\varphi\|\varphi\|^2$  for  $\underline{x}^*$  balanced.  $\square$

**Lemma B.5** *Let  $(\underline{x}^*, F)$  be an solution/instance pair distributed according to the planted model  $\text{pCSP}(n, \alpha, p)$  and assume*

$$\frac{\tilde{\Omega}_k}{k} (\log k)(1 + \varepsilon) \leq \alpha \leq \Omega_k (\log 2)(1 - \varepsilon). \quad (57)$$

*Then there exists constants  $0 < \theta_1 < \theta_2 < 1$ , and  $c, c' > 0$  such that, with probability at least  $1 - e^{-cn}$  the following happens. The instance  $F$  does not admit any solution  $\underline{x}$  with  $n\theta_1 \leq \underline{x} \cdot \underline{x}^* \leq n\theta_2$ , and the number of solutions with  $\underline{x} \cdot \underline{x}^* \geq n\theta_2$  is at most  $e^{-nc'} \mathbb{E}Z(F)$  (expectation is here with respect to the uniform model).*

**Proof.** In view of Lemma B.4 it is sufficient to show that there  $\theta_* \in (0, 1)$  such that:

- (a)  $\Psi(\theta_*) < 0$ .
- (b)  $\sup_{\theta \in [\theta_*, 1]} \Psi(\theta) < \log 2 + \alpha \log \mathbb{E}\|\varphi\|^2$ .

In order to prove (a), we first notice that, for any  $\varepsilon \in (0, 1/2)$

$$\frac{\mathbb{E}_\varphi(\varphi, T_\theta \varphi)}{\mathbb{E}_\varphi\|\varphi\|^2} \leq 1 - \frac{1}{(1 + \varepsilon)\tilde{\Omega}_k} + \frac{1}{(1 + \varepsilon)\tilde{\Omega}_k} e^{-k(1+\varepsilon)(1-\theta)}, \quad (58)$$

provided  $\theta > 1 - \varepsilon$ . Indeed both sides equal 1 at  $\theta = 1$ . Further, the derivative of the left hand side can be estimated as

$$\begin{aligned} \frac{d}{d\theta} \frac{\mathbb{E}_\varphi(\varphi, T_\theta \varphi)}{\mathbb{E}_\varphi\|\varphi\|^2} &= 2k \frac{\mathbb{E}_\varphi(\varphi^{(1)}, T_\theta \varphi^{(1)})}{\mathbb{E}_\varphi\|\varphi\|^2} \geq 2k \frac{e^{-k(1+\varepsilon)(1-\theta)} \mathbb{E}_\varphi\|\varphi^{(1)}\|^2}{\mathbb{E}_\varphi\|\varphi\|^2} \\ &= k e^{-k(1+\varepsilon)(1-\theta)} \frac{2\mathbb{E}_\varphi I_1(\varphi)}{1 - 2\mathbb{E}_\varphi I_1(\varphi)} \geq \frac{d}{d\theta} \left\{ 1 - \frac{1}{(1 + \varepsilon)\tilde{\Omega}_k} + \frac{1}{(1 + \varepsilon)\tilde{\Omega}_k} e^{-k(1+\varepsilon)(1-\theta)} \right\}. \end{aligned}$$

Here we used the following inequality, valid for any  $f : \{+1, -1\}^k \rightarrow \{0, 1\}$ , provided  $\theta > 1 - \varepsilon$

$$(f, T_\theta f) = \sum_Q |f_Q|^2 \theta^{|Q|} \geq \|f\|^2 \theta^k \geq \|f\|^2 e^{-k(1+\varepsilon)(1-\theta)}. \quad (59)$$

Let  $\alpha = (1 + \varepsilon)(\tilde{\Omega}_k/k)\gamma \log k/k$ , and  $\theta_* = 1 - \omega_*/k$ . Equation (58) implies

$$\Psi(\theta_* = 1 - \omega_*/k) \leq H(\omega_*/k) - \frac{\gamma}{k} (\log k) + \frac{\gamma}{k} (\log k) e^{-(1+\varepsilon)\omega_*}, \quad (60)$$

for all  $\varepsilon > \omega_*/k$ . If we fix  $\varepsilon = \omega_{\max}/k$ , and let  $k \rightarrow \infty$ , we finally obtain (for  $\omega \in (0, \omega_{\max})$ )

$$\Psi(\theta_* = 1 - \omega_*/k) \leq \left\{ \omega_* - \gamma + \gamma e^{-\omega_*} \right\} \frac{\log k}{k} + O(k^{-1}). \quad (61)$$

As soon as  $\gamma > 1$ , we can find  $\omega_*$  such that  $\omega_* - \gamma + \gamma e^{-\omega_*} < 1$  (just take  $\omega_* = \log \gamma$ ). Further,  $\sup_{\theta \in [\theta_*, 1]} \Psi(\theta) = O(1/k)$  which is smaller than  $\log 2 + \alpha \log \mathbb{E} \|\varphi\|^2$  for  $k$  large enough and  $\alpha < \Omega_k(\log 2)(1 - \varepsilon)$ .  $\square$

**Proof (Theorem 3.3).** Consider a random instance from the  $\text{CSP}(n, \alpha, p)$  ensemble, and sample a solution  $\underline{x}^*$  uniformly at random. By Lemma B.5 and Theorem B.3, with high probability there is no solution  $\underline{x}$  such that  $\underline{x} \cdot \underline{x}^* \in [n\theta_1, n\theta_2]$ . Declare the cluster of  $\underline{x}^*$ ,  $\mathcal{C}(\underline{x}^*)$  to be the set of solutions  $\underline{x}$  such that  $\underline{x} \cdot \underline{x}^* \geq n\theta_2$ . It will contain an exponential small fraction of solutions.

The same operation can be repeated  $e^{n\delta}$  times. Since each cluster thus constructed is exponentially small, for  $\delta$  small enough the probability that any of the two clusters intersects is exponentially small.  $\square$

## C Sharp Threshold Results for CSPs

Recall that in the previous section, we appealed crucially in two places to certain sharp transition behavior of the CSP's under consideration. We furnish the requisite references and details here.

Since we are interested in the behavior of binary  $k$ -CSP's for large  $k$ , in what follows we may safely assume that  $k \geq 3$ . Once again for simplicity, let  $F = F_k(n, \alpha n)$  denote a random binary CSP( $n, \alpha, p$ ) on  $n$  variables and  $\alpha n$  clauses, and the distribution  $p$  over clauses satisfying the main conditions (1)–(4) mentioned in Section 3. As is customary, for the SAT-UNSAT threshold to be meaningful, we also assume that  $p$  satisfies the following elementary condition.

**5. Unsatisfiability of the ensemble.** For every  $\epsilon = \pm 1$ , there is at least one clause  $g$  with  $p_g > 0$  such that  $g(\epsilon, \dots, \epsilon) = 0$ . (Note that by the balance condition (2), necessarily  $g(-\epsilon, \dots, -\epsilon) = 0$ ).

Building on their previous work, Creignou and Daude recently showed [CD09] that the satisfiability of  $F_k(n, \alpha n)$  undergoes a sharp transition, except when the formula contains a function of one of the following two types:

(i) A Boolean function  $f$  *strongly depends on one component* if there exist  $\epsilon \in \{+1, -1\}$  and  $i$  with  $1 \leq i \leq k$  such that  $(x_1, \dots, x_n) \in \{+1, -1\}^n$  and  $f(x_1, \dots, x_n) = 1$  imply that  $x_i = \epsilon$ .

(ii) A Boolean function  $f$  *strongly depends on a 2-XOR-relation* if there exist  $i, j$  with  $1 \leq i \neq j \leq k$  such that  $(x_1, \dots, x_n) \in \{+1, -1\}^n$  and  $f(x_1, \dots, x_n) = 1$  imply that  $x_i \oplus x_j = 1$ .

**Theorem C.1 (CD09)** *With  $F = F_k(n, \alpha n)$  and  $p$  satisfying (5) above, the transition from SAT( $F$ ) to UNSAT( $F$ ) is sharp if and only if  $F$  contains no function strongly depending on one component and no function strongly dependent on a 2-XOR-relation.*

Note that we had used this result in completing the proof of the lower bound in Proposition 3.1, in Appendix A.

We now furnish various details needed to justify that the property of having an exponential number of solutions has a sharp threshold. Recall that this was needed to boost the Lemma B.2 (see Appendix B) in

the proof of the clustering threshold, to show that the probability once bounded away from 0, is actually tending to 1, as the problem size  $n$  went to infinity.

Let  $\Phi$  be a formula on the variables  $y_1, \dots, y_l$  that can be constructed from our ensemble, let  $X = \{x_1, \dots, x_n\}$  be a set of  $n$  variables (disjoint from  $\{y_1, \dots, y_l\}$ ), and let  $\Phi_n$  denote the set of all formulas that result after substituting  $l$  distinct variables from  $X$  and replacing them in  $\Phi$ . Given a CSP ensemble  $F$  on  $n$  variables, let  $F \oplus \Phi$  be equal to  $F \wedge \Phi^*$ , where  $\Phi^*$  is a random formula chosen uniformly from  $\Phi_n$ .

We say a random ensemble  $F$  has the property  $\mathcal{A}_B = \mathcal{A}_B(F)$  if  $F$  has fewer than  $\frac{1}{2}B^n$  satisfying assignments. We want to prove the following:

**Lemma C.2** *For any  $B \in [1, 2)$  there is a sequence  $t_n^B$  such that for any  $\epsilon > 0$ ,*

$$\begin{aligned} \lim_{n \rightarrow +\infty} \mathbf{P}(F_k(n, (1 - \epsilon)t_n^B) \text{ has property } \mathcal{A}_B) = 0, \quad \text{and} \\ \lim_{n \rightarrow +\infty} \mathbf{P}(F_k(n, (1 + \epsilon)t_n^B) \text{ has property } \mathcal{A}_B) = 1. \end{aligned} \tag{62}$$

Note that  $\mathcal{A}_B$  is a monotone property, since whenever  $F$  has the property, then  $F \wedge F'$  will have the property for any formula  $F'$  on the variables  $\{x_1, \dots, x_n\}$ . We will use the following theorem of Friedgut [F05] to prove that  $\mathcal{A}_B$  has a ‘sharp threshold’, in the sense of Lemma C.2.

**Theorem C.3 (F05)** *Suppose that  $\mathcal{A}_B$  does not have a sharp threshold. Then, there exists  $\alpha > 0$ , a formula  $\Phi$ , and for any  $n_0 > 0$ , there exist  $n > n_0$ ,  $m > 0$ , and a formula  $F$  with variables  $x_1, \dots, x_n$  such that all of the following hold:*

*T1 .  $\mathbf{P}(F \oplus \Phi \text{ has the property } \mathcal{A}_B) > 1 - \alpha$ .*

*T2 .  $\alpha < \mathbf{P}(F_k(n, m) \text{ has the property } \mathcal{A}_B) < 1 - 3\alpha$ .*

*T3 . With probability at least  $\alpha$ , a random formula  $F_k(n, m)$  contains an element of  $\Phi_n$  as a subformula.*

*T4 .  $\mathbf{P}(F \wedge F_k(n, 2 \log n) \text{ has the property } \mathcal{A}_B) < 1 - 2\alpha$ .*

A first observation is the subtle fact that Theorem C.3 is originally stated in terms of a parametric Bernoulli model, while our model is Binomial. But it is the case, by standard arguments, that we can translate results concerning the existence of a sharp threshold of monotone properties from one model to other, provided that  $m$  is of order  $\Omega(n)$ . We will prove that this is the case, in step (1) below.

An important fact that we will use throughout is that, because of the feasibility condition, a *pure literal* reduction scheme exists: Suppose that  $x_l$  is a variable that appears *only once* in a formula  $F = C_1 \wedge \dots \wedge C_m$ , say, in the clause  $C_1 = f(x_l, x_{i_1}, \dots, x_{i_{k-1}})$ . Then, any satisfying assignment  $\chi : [n] \setminus \{l\} \rightarrow \{\pm 1\}$  of  $C_2 \wedge \dots \wedge C_m$  can be extended to a satisfying assignment  $\bar{\chi} : [n] \rightarrow \{\pm 1\}$  of  $C_1 \wedge C_2 \wedge \dots \wedge C_m$ , by setting  $\bar{\chi}(l)$  to the appropriate value (due to feasibility), such that  $f(\bar{\chi}(l), \chi(i_1), \dots, \chi(i_{k-1})) = 1$ .

Notice that using iteratively a pure literal reduction scheme, we can find a satisfying assignment for the formula if we can iteratively find a variable contained once in the formula, eliminate the clause containing the variable and proceed again with the new formula, until obtaining an empty formula. This procedure is equivalent to such of finding the 2-core of the associated hypergraph [M05], and in fact, it is the case that if the associated hypergraph has an empty 2-core, then this pure literal reduction scheme will be successful in finding a satisfying assignment.

The approach we will use to prove Lemma C.2 goes along the lines of [AC08, Lemma 13], with some variations that follow the work of Creignou and Daude in [CD02], [CD04] and [CD09]. As is standard in these proofs, in the sequel we will assume the existence of  $\alpha$ ,  $\Phi$ ,  $n$  and  $m$  satisfying **T1-T3** and to conclude

that the property  $\mathcal{A}_B$  has a sharp threshold, we will prove that **T4** cannot hold. Notice that we can always assume that  $n$  is large enough, by choosing  $n_0$  appropriately. We will divide the core of the proof in three steps: In the first one we determine the correct scaling of  $m$ . In the second step we prove that the small formula  $\Phi$  is indeed satisfiable. And, in the last step, we proceed in concluding that **T4** does not hold, completing the contradiction argument.

**(1) Scaling of  $m$ :**

*Lower bound:* Notice that for  $m \equiv \epsilon n/k$ , necessarily  $(1 - \epsilon)n$  variables do not appear in  $F_k(n, m)$ , so that if  $F_k(n, m)$  is satisfiable, it contains at least  $2^{(1-\epsilon)n}$  satisfying assignments. But, following [M05], there is a constant  $c^*$  such that if  $m < c^*n$ , then the hypergraph associated to  $F_k(n, m)$  w.h.p. does not have a 2-core, and as mentioned before, the pure literal reduction is successful in finding a satisfying assignment. This proves, by choosing  $\epsilon$  small enough, that for  $m \equiv \epsilon n/k$ , w.h.p.,  $F_k(n, m)$  has at least  $2^{(1-\epsilon)n} \geq \frac{1}{2}B^n$  satisfying assignments. Therefore, by **T2**, it should be the case that  $m = \Omega(n)$ .

*Upper bound:* From the first moment estimates in the present paper, we have that there is a constant  $C_p$  (depending only on  $p$ ), such that w.h.p, a random formula  $F_k(n, C_p n)$  is not satisfiable. Therefore (by **T2**), due to the monotonicity of  $\mathcal{A}_B$  it should be the case that  $m = O(n)$ .

**(2) Satisfiability of  $\Phi$ :** Given a formula  $\Phi$ , define  $\mathbf{v}(\Phi)$  to be the number of variables in  $\Phi$ , and  $\mathbf{w}(\Phi)$  to be the number of clauses in  $\Phi$ . By an easy counting, for any  $t \geq 1$ , if  $m = O(n)$ , then the probability that a random formula  $F_k(n, m)$  contains a subformula  $\Phi$  with  $w(\Phi) \leq t$  and such that  $\mathbf{v}(\Phi) \leq (k - 1)\mathbf{w}(\Phi) - 1$ , goes to zero as  $n \rightarrow +\infty$ . Now, if  $\Phi$  is unsatisfiable, then it contains a minimal unsatisfiable formula  $\psi$  with  $\mathbf{w}(\psi) \leq t$ , and therefore, by the previous conclusion, by **T2** and **T3**, we have that  $\mathbf{v}(\psi) > (k - 1)\mathbf{w}(\psi)$  w.h.p. Then, using [CD02, Lemma 5.2],  $\psi$  has either a constraint with  $k - 1$  variables appearing only once, or it is unicyclic. In either case, for  $k \geq 3$ , there is at least one variable appearing only once in the formula, therefore the pure literal reduction operates, contradicting the minimality of  $\psi$ .

**(3) Contradicting T4:**

*Step 3a:* By **T3** and the conclusion of step (1),  $\Phi$  is w.h.p. satisfiable. Let  $\{y_1, \dots, y_l\}$  be the variables appearing in  $\Phi$ , and let  $\sigma : \{1, \dots, l\} \rightarrow \{\pm 1\}$  be a fixed satisfying assignment of  $\Phi$ . We say that a satisfying assignment  $\chi$  of  $F$  is compatible with a tuple  $(z_1, \dots, z_l) \in [n]^l$  if  $\chi(z_i) = \sigma(i)$  for all  $i = 1, \dots, l$ . Furthermore, we say that the tuple  $(z_1, \dots, z_k)$  is *bad* if  $F$  has fewer than  $\frac{1}{2}B^n$  satisfying assignments compatible with  $(z_1, \dots, z_l)$ . Notice that by **T1**, there are at least  $(1 - \alpha)n^l$  bad tuples.

*Step 3b:* By the Erdos-Simonovits theorem [ES82], if  $l$   $k$ -tuples  $(w_1^1, \dots, w_1^k), \dots, (w_l^1, \dots, w_l^k)$  are chosen uniformly at random and independently from  $n^k$ , then with probability at least  $\gamma'$ , for every function  $f : [l] \rightarrow [k]$ , the tuple  $(w_1^{f(1)}, \dots, w_l^{f(l)})$  is a bad tuple. In particular we have that with probability at most  $(1 - p_g^l \gamma')^{(\log n)/l}$ , a random formula  $F_k(n, \log n)$  will not contain  $l$  clauses  $C_1, \dots, C_l$  satisfying

- (i)  $C_i = g(v_i^1, \dots, v_i^k)$  for  $i = 1, \dots, l$ , where  $g$  is the boolean function whose existence is implied by condition **5**.
- (ii) For every function  $f : [l] \rightarrow [k]$ , the  $l$ -tuple  $(v_1^{f(1)}, \dots, v_l^{f(l)})$  is bad.

Therefore, by choosing  $n$  large enough, the probability that a random formula  $F_k(n, \log n)$  contains clauses satisfying (i) and (ii) is at least  $1 - \alpha$ .

*Step 3c:* Let  $C_1, \dots, C_l$  be clauses satisfying **(i)** and **(ii)**, and let  $\chi : [n] \rightarrow \{\pm 1\}$  be a satisfying assignment of  $F \wedge C_1 \wedge \dots \wedge C_l$ . Then note that for every  $i = 1, \dots, l$ , there exists an  $f(i)$  such that  $\chi(v_i^{f(i)}) = \sigma(i)$ . Otherwise, for some  $i$ , and all  $j = 1, \dots, k$ ,  $\chi(v_i^j) = -\sigma(i)$ , which implies that  $\chi$  does not satisfy  $C_i$ , which is a contradiction. It now follows that  $\chi$  is compatible with  $(v_1^{f(1)}, \dots, v_l^{f(l)})$ . Therefore, we conclude that every satisfying assignment of  $F \wedge C_1 \wedge \dots \wedge C_l$  is compatible with  $(v_1^{f(1)}, \dots, v_l^{f(l)})$  for some function  $f : [l] \rightarrow [k]$ . But, by condition **(ii)**, every one of these  $l$ -tuples is bad, and therefore, each one does not have more than  $\frac{1}{2}B^n$  satisfying assignments compatible with them. As a result,  $F \wedge C_1 \wedge \dots \wedge C_l$  does not have more than  $\frac{1}{2}k^l B^n$  satisfying assignments. Moreover, combining Step 2b and Step 2c, we conclude that with probability at least  $1 - \alpha$ ,  $F \wedge F^*$  contains at most  $\frac{1}{2}k^l B^n$  satisfying assignments, where  $F^*$  is a random  $F_k(n, \log n)$  formula.

*Step 3d:* Given a satisfying assignment  $\chi : [n] \rightarrow \{\pm 1\}$ , with probability at least  $2^{1-k}$ , the clause  $g(v_1, \dots, v_k)$ , where  $(v_1, \dots, v_k)$  is chosen uniformly are random from  $[n]^k$ , will not be satisfied by  $\chi$ . In particular a random clause will be satisfied by  $\chi$  with probability at most  $1 - p_g 2^{1-k}$ . More generally, a random  $F_k(n, \log n)$  will be satisfied by  $\chi$  with probability at most  $(1 - p_g 2^{1-k})^{\log n} \leq \frac{1}{n^{c_k}}$ , where  $c_k = p_g 2^{1-k}$ . Therefore, if  $F^{**}$  is a  $F_k(n, \log n)$  random formula independent of  $F^*$ , we have that

$$\mathbf{E} \left[ \# \text{sat. assign. of } F \wedge F^* \wedge F^{**} \mid \# \text{ sat. assign. of } F \wedge F^* \leq \frac{1}{2} k^l B^n \right] \leq \frac{1}{2n^{c_k}} k^l B^n,$$

and therefore, by Markov's inequality

$$\mathbf{P} \left[ \# \text{sat. assign. of } F \wedge F^* \wedge F^{**} \geq \frac{1}{2} B^n \mid \# \text{ sat. assign. of } F \wedge F^* \leq \frac{1}{2} k^l B^n \right] \leq \frac{k^l}{n^{c_k}},$$

which is less than  $\alpha/2$  for  $n$  large enough. Thus, combining the conclusion of Step 2c and the previous formula, we obtain

$$\mathbf{P} \left[ \# \text{sat. assign. of } F \wedge F^* \wedge F^{**} \geq \frac{1}{2} B^n \right] \geq 3\alpha/2,$$

and this contradicts **T4**, thereby proving that property  $\mathcal{A}_B$  has a sharp threshold. □