

## QUIET PLANTING IN THE LOCKED CONSTRAINT SATISFACTION PROBLEMS

LENKA ZDEBOROVÁ\* AND FLORENT KRZAKALA†

**Abstract.** We study the planted ensemble of locked constraint satisfaction problems. We describe the connection between the random and planted ensembles. The use of the cavity method is combined with arguments from reconstruction on trees and the first and the second moment considerations. Our main result is the location of the hard region in the planted ensemble. In a part of that hard region instances have with high probability a single satisfying assignment.

**Key words.** Constraint Satisfaction Problems, Planted Random Ensemble, Belief Propagation, Reconstruction on Trees, Instances with a Unique Satisfying Assignment.

**AMS subject classifications.** 90C27 68Q25 05C80

Constraint Satisfaction Problems (CSPs) are very general in their nature: Consider a set of  $N$  discrete variables and a set of  $M$  Boolean constraints; the problem consists in finding a configuration of variables that satisfies all the constraints or in proving that no such configuration exists. As such, CSPs are a subject of interest in many different fields such as computer science, discrete mathematics, physics, engineering and computational biology. Random ensembles of CSPs are a fertile source of research activity; as hard benchmarks they serve for testing new algorithmic ideas [5, 31], they are used to create efficient coding schemes [14, 15], to model complex glass forming liquids [4, 23], or to understand the origin of average computational hardness [30, 43]. Combining know-how from many branches of mathematics, computer science and statistical physics seems to be a fruitful strategy for the understanding of these stunning objects and of their very rich behavior.

The most commonly studied random ensembles of CSPs are created by choosing the graph of variables and constraints as a random bipartite graph with a certain left and right degree distribution. Another natural way of creating a random instance, called planting, is to first assign a configuration to all variables, and then to choose only constraints compatible with this configuration. Both these ensembles can be useful to mimic instances created in some practical applications, as e.g. the low density parity check codes [14]. In particular planted instances may be created in adaptive situations when only constraints satisfied by the current state of variables can be added.

By planting, we create by definition a satisfiable instance. Such instances are particularly useful as benchmarks to evaluate the performance of incomplete solvers, such as stochastic local search [38]. Based on the example of the planted  $K$ -satisfiability problem it is, however, often anticipated that the planted ensemble is algorithmically easier than the random one because a bias towards the planted assignment is created in the graph. Also, for most of the studied problems, it was proven that at large density of constraints is it indeed easy to find a satisfying assignment near to the planted one, see e.g. [3, 7, 13]. On the other hand if the planted ensemble would be algorithmically hard in some region of parameters than these instances could serve

---

\*Theoretical Division and Center for Nonlinear Studies, Los Alamos National Laboratory, NM 87545 USA ([lenka@lanl.org](mailto:lenka@lanl.org)).

†CNRS and ESPCI ParisTech, 10 rue Vauquelin, UMR 7083 Gulliver, Paris 75000 France ([fk@espci.fr](mailto:fk@espci.fr)). Theoretical Division and Center for Nonlinear Studies, Los Alamos National Laboratory, NM 87545 USA.

as one-way functions and have applications in cryptography. Yet, compared to the random ensemble, relatively little is known about the existence, size and properties of algorithmically hard regions in the planted ensemble.

In this paper we study a way of planting an assignment which changes only in a minimal way the properties of the random ensemble. We call this a *quiet planting*. Although the concept of *quiet planting* was introduced in [24], many results were actually first demonstrated and used as a tool for proofs in [2]. Both these works, however, were mainly concentrated on the coloring problem (and on hyper-graph bi-coloring). In the present work we generalize this idea and we will focus on quiet planting in the so-called locked CSPs, introduced recently in [43, 42]. On one hand, the locked CSPs have a very interesting phase diagram which description-wise is much simpler than the one of the graph coloring or K-satisfiability. On the other hand they are much harder algorithmically and the boundaries between the easy and hard regions are, unlike in the coloring or K-satisfiability, relatively well understood (at least on the heuristic level of the cavity method [43, 42]). This special behavior stems from the fact that in the locked problems the space of solutions consists of points separated by an extensive (i.e.  $\Theta(N)$ ) Hamming distance instead of clusters of solutions.

Here, we combine the idea of quiet planting with the special behavior of the locked CSPs and obtain random CSPs ensembles with many interesting properties that are summarized in Sec. 1 in the context of related works. In Sec. 2 we introduce the necessary definitions and notations, and in Sec. 3 we summarize the phase diagram of the locked problems derived in [43, 42]. In Sec. 4 we argue about the equivalence between the random and the planted ensemble based on heuristic considerations and on a second moment computation. In Sec. 5 we describe the interesting phase where instances of our problems have with large probability a single solution. Finally, in Sec. 6 we discuss the algorithmic hardness of the planted instances, and compute the critical degree beyond which instances become easy to solve. We conclude the paper with a list of open problems in Sec. 7.

**1. Main results and related works.** The results of this paper apply to the *factorized locked* CSPs, see Defs. 2.2, 2.4. We list in six points the most important contributions of the present article:

- (i) **The planted configuration is equivalent to randomly sampled satisfying assignment:** The idea of *quiet planting* is to plant a configuration that in many important properties does not differ from a satisfying assignment sampled uniformly at random on the resulting graph. Such a problem is closely related to the reconstruction on trees [10, 34, 26] where we consider an assignment taken uniformly at random from all the satisfying ones. Constructing such an assignment on a tree is always possible in polynomial time, as the exact marginal probabilities can be obtained via the belief propagation (BP) algorithm. On random graphs, uniform sampling is in general exponentially costly. However, the quiet planting can be achieved asymptotically on the factorized CSPs, see Def. 2.4. In a statistical physics language quiet planting is possible whenever the quenched entropy equals the annealed one, see condition (4.2). The possibility of quiet planting and its relation to a concentration property (4.2) was previously discussed for other constraint satisfaction problems in [32, 2, 24].
- (ii) **Equivalence between the planted and the random ensemble in the satisfiable phase:** Many properties of the planted ensemble created via quiet planting can be deduced from the properties of the purely random ensemble.

Among others, in the satisfiable phase the random and planted ensembles are asymptotically equivalent, see Def. 4.2. Such equivalence can also be established rigorously based on a second moment argument, as in [2]. In the factorized locked problems treated in this article the second moment is able to pin the satisfiability threshold sharply and hence the equivalence between the planted and the random ensemble holds in the whole satisfiable phase.

- (iii) **Equivalence between planted and satisfiable ensembles:** Based on heuristic (cavity) arguments we conjecture that the planted ensemble is in the factorized locked problems asymptotically equivalent to the ensemble of satisfiable instances in the whole range of  $\Theta(1)$  constraint densities. In particular, this means that one can recognize easily almost all the rare satisfiable instances above the robust reconstruction threshold by using belief propagation. We stress here that we do not expect this equivalence to be true in the non-locked CSPs, such as graph coloring or K-satisfiability.
- (iv) **Easy-hard-easy transition:** Next to the interesting conceptual results, the most important practical result is the identification and the location of the region where the instances from the planted ensemble of the factorized locked problems are on average computationally hard. We show that an easy-hard-easy pattern for finding a solution appears in the planted ensemble as the constraint density is increased, where *easy* means that an on-average polynomial algorithm is known, while *hard* means that no on-average polynomial algorithm is known (and maybe does not exist). We conjecture that the two boundaries of the hard phase correspond to two different reconstruction thresholds – the onset of hardness coincides with the *small noise* reconstruction threshold [42], called the dynamical transition in the physics literature [22], and the end of the hard region is given by the threshold for the *robust* reconstruction [16]. This last point also corresponds to the Kesten-Stigum bound for the canonical reconstruction on trees [19, 20], and to the spin glass local instability in the purely random ensemble [29]. We also show that outside the hard region algorithms based on belief propagation are able to find solutions efficiently. In particular in the high average degree easy region, the belief propagation algorithm converges directly to the planted solution.
- (v) **Hard satisfiable benchmarks:** Given we have located the values of parameters where the instances of the planted ensemble are hard, these can serve as very challenging satisfiable benchmarks. Such benchmarks are in particular interesting for the evaluation of regions in which incomplete solvers work—or do not work—in polynomial (linear) time. Note that for complete exhaustive solvers the locked problems are not necessarily harder than the canonical K-satisfiability. On the contrary, locked constraints produce more implications when a variable is fixed, hence exhaustive branch-and-bound techniques might come to a decision relatively faster than in the random K-satisfiability.
- (vi) **Instances with unique satisfying assignment (USA):** We show that beyond the threshold corresponding to the satisfiability threshold in the random ensemble the planted instances have with high probability a single satisfying assignment (or a pair of them in case a global symmetry is present). Moreover depending on the constraint density these USA instances can be found in the hard or in the easy region. Some USA instances are extensively used in evaluation of quantum algorithms, see e.g. [40, 12]. In these previous works these instances are, however, generated with exponential cost, and their classical

typical computational hardness has not been evaluated.

A large part of our results is based on the heuristic cavity method approach [28]. We were also able to prove part of our results for the  $R$ -in- $K$  SAT problem on random regular graphs using computations of the second moment and the expander property. This includes some results about the equivalence between the planted and random ensembles in the satisfiable phase, and the uniqueness of the satisfying assignment in the unsatisfiable phase. Completing and extending these proofs to the other locked factorized CSPs should be possible although more involved.

TABLE 1.1

*Sketchy summary of the properties of the different phases in the random ensemble of the factorized locked problems, the parameter  $l$  is the average number of constraints in which a variable appears. The three thresholds  $l_d$ ,  $l_s$  and  $l_l$  are defined in detail later in the paper.*

RANDOM	$l < l_d$	$l_d < l < l_s$	$l_s < l < l_l$	$l_l < l$
BP, random init.	converges	converges	converges	does not
BP fixed point	uniform	uniform	uniform	×
# of solutions	exponential	exponential	none	none
finding solution	easy	hard	×	×
reconstruction	not possible	possible	×	×

TABLE 1.2

*The same as Tab. 1.1 for the random planted ensemble, its definition is given in Sec. 4.*

PLANTED	$l < l_d$	$l_d < l < l_s$	$l_s < l < l_l$	$l_l < l$
BP, random init.	converges	converges	converges	converges
BP fixed point	uniform	uniform	uniform	planted
BP, planted init.	converges	converges	converges	converges
BP fixed point	uniform	planted	planted	planted
# of solutions	exponential	exponential	one/two	one/two
finding solution	easy	hard	hard	easy
reconstruction	not possible	possible	possible	possible
robust recons.	not possible	not possible	not possible	possible

**2. Definitions and notations.** In this section we specify the class of constraint satisfaction problems to which our results apply. The crucial notions will be the definition of a *locked* [43] and *factorized* constraint satisfaction problem. It is only on the factorized problems where there is a very close relation between the usual random and the planted ensemble, as discussed in [24]. It is also the fact that in the locked problems solutions (i.e. satisfying assignments) are mutually far from each other in terms of their Hamming distance [43] that makes them particularly interesting for considerations in this context.

**DEFINITION 2.1.** *A constraint  $a$  containing  $K$  variables, the domain of each variable being  $X$ , is a function from  $X^K$  to  $\{0, 1\}$ . If the function evaluates to 1 (0) we say that constraint  $a$  is satisfied (not satisfied). A constraint is locked if and only if there are no two satisfying assignments of variables which would differ in a single value (out of the  $K$  ones). In this paper we will consider for concreteness binary variables, that is  $X = \{0, 1\}$ , but the results are generalizable to larger domain sizes.*

DEFINITION 2.2. A constraint satisfaction problem consists in deciding if there exists a configuration of  $N$  variables which satisfies simultaneously a set of  $M$  constraints. A constraint satisfaction problem is called locked if and only if all the  $M$  constraint are locked and each of the  $N$  variables belongs to at least two different constraints. Thus anytime we speak about a locked problem we implicitly suppose that the corresponding factor-graph does not have any leaves (variables of degree one). The degree of a variable is defined as the number of constraints to which the variable belongs, while the degree of a constraint is the number of variables it contains.

We shall illustrate our findings on the so called occupation constraint satisfaction problems [33, 42].

DEFINITION 2.3. In occupation problems every constraint  $a$  depends only on the sum of the variables it contains. Thus every occupation constraint containing  $K_a$  variables can be characterized by a binary  $K_a + 1$  component vector  $A_a$  such that the constraint is satisfied if and only if the sum  $r$  of the  $K_a$  variables is such that  $A_a(r) = 1$ .

An occupation constraint  $a$  is locked if and only if for all  $i = 0, \dots, K_a - 1$  we have  $A_a(i)A_a(i+1) = 0$ . We will consider occupation problems where every constraint contains  $K$  variables and is given by the same vector  $A$ . To give an example of this notation, the vector  $A = 0100$  corresponds to the 1-in-3 SAT problem (also called exact cover), which is indeed locked. The vector  $A = 0110$  corresponds to the hypergraph bi-coloring problem, which is not locked (since there are two neighboring 1s). Many other examples can be found in [42]. For problems which do not have other name established in the literature, we will use the notation  $i$ -or- $j$ -...-in- $K$  SAT for a vector  $A$  with non-zero components  $A(i), A(j)$ , etc.

Let us now write the belief propagation (BP) equations [35, 25, 27] for the occupation constraint satisfaction problems. The basic quantities in BP are messages. We define  $\mu_{s_i}^{a \rightarrow i}$  as the probability (over all satisfying assignments) that the variable  $i$  has value  $s_i$  given that  $i$  belong only to constraint  $a$ . The BP equations approximate these probabilities  $\mu_{s_i}^{a \rightarrow i}$  by messages  $\psi_{s_i}^{a \rightarrow i}$  by assuming that the factor graph [25] underlying the CSP is a tree

$$\psi_{s_i}^{a \rightarrow i} = \frac{1}{Z^{a \rightarrow i}} \sum_{\{s_j\}} \delta_{A(s_i + \sum_j s_j), 1} \prod_{j \in \partial a - i} \prod_{b \in \partial j - a} \psi_{s_j}^{b \rightarrow j}, \quad (2.1)$$

where  $Z^{a \rightarrow i}$  is a normalization constant assuring  $\psi_1^{a \rightarrow i} + \psi_0^{a \rightarrow i} = 1$ ,  $\partial a$  is the set of neighbors of  $a$ ,  $\partial a - i$  are neighbors of  $a$  except  $i$ , and the sum over  $\{s_j\}$  is over all values variables  $s_j$  can take. Fig. 2.1 shows the corresponding part of the factor graph.

We define  $\nu^i$  to be the probability (over all satisfying assignments) that a variable  $i$  is occupied. The BP estimate of the probability that a variables  $i$  is occupied is

$$\chi^i = \frac{\prod_{a \in \partial i} \psi_1^{a \rightarrow i}}{\prod_{a \in \partial i} \psi_1^{a \rightarrow i} + \prod_{a \in \partial i} \psi_0^{a \rightarrow i}}, \quad (2.2)$$

Note that if an assignment  $\{\sigma\}$  is a solution of the locked problem, then  $\psi_{\sigma_i}^{a \rightarrow i} = 1$ ,  $\psi_{-\sigma_i}^{a \rightarrow i} = 0$  is a fixed point of the BP equations (2.1). If the underlying factor-graph is a tree then the fixed point of the BP equations is unique and  $\mu_{s_i}^{a \rightarrow i} = \psi_{s_i}^{a \rightarrow i}$  and  $\nu^i = \chi^i$  in the fixed point (note that on a tree the problem is not locked). On a graph with cycles this is not the case in general. We will call a BP fixed point asymptotically ( $N \rightarrow \infty$ ) exact on a random ensemble of graphs if  $\mu_{s_i}^{a \rightarrow i} = \psi_{s_i}^{a \rightarrow i} + o(1)$  and  $\nu^i = \chi^i + o(1)$  for almost all  $a$  and  $i$  with high probability, where  $N$  is the number of nodes in the graph.

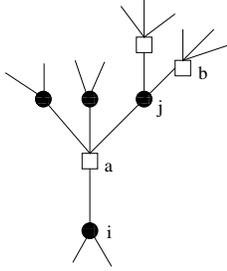


FIG. 2.1. Part of the factor graph to illustrate the meaning of indices in the belief propagation equations (2.1).

All our results are restricted to ensembles of constraint satisfaction problems where at least one BP fixed point is asymptotically exact. We define such ensembles in the remaining of this section. A graph ensemble is locally tree-like if the shortest loop going through a random node has w.h.p. length diverging as  $N \rightarrow \infty$ . Families of sparse random graphs, i.e. the degree distribution of variables  $Q(l)$  does not depend on  $N$ , are locally tree-like as long as the mean of  $Q(l)$  is finite. The variable degree distributions we will be using mostly are:

- Regular  $Q(l) = \delta_{L,l}$ .
- Truncated Poisson  $Q(0) = Q(1) = 0$ ,  $Q(l) = c^l / [(e^c - 1 - c)l!]$  for  $l \geq 3$ . The average degree in this case is  $\bar{l} = c(1 - e^{-c}) / [1 - (1 + c)e^{-c}]$ .

In order to generate random graphs with a given variable degree distribution, one can apply the following algorithm:

- Repeat until  $KM = \sum_{i=1}^N l_i$ : Draw  $N$  random numbers  $l_i$  from distribution  $Q(l)$ .
- Consider  $K$  legs going from every constraint, order them arbitrarily and index them from 1 to  $KM$ , consider  $l_i$  legs from every variable  $i$ , order them arbitrarily and index them from 1 to  $KM$ .
- Repeat until there are no double edges: Draw a random permutation  $\pi$  of  $KM$  numbers and connect  $i$ -th leg from constraints with  $\pi(i)$ -th leg from variables.

We define an iteration of the belief propagation algorithm as taking all the  $KM$  edges  $ai$  in a random order and updating the message  $\psi_{s_i}^{a \rightarrow i}$  according to eq. (2.1).

DEFINITION 2.4. A given instance of a constraint satisfaction problem is factorized if and only if the belief propagation equations initialized randomly converge almost surely (with probability approaching one as the number of variables  $N \rightarrow \infty$ ) to a uniform fixed point, i.e., the value of  $\psi^{a \rightarrow i}$  is the same for almost all edges  $ai$ .

Note that it is a non-trivial task to provably decide if a problem satisfies this definition, and the answer depends on the degree distribution. In practice we generate a large random instance of the problem, initialize BP randomly and iterate. We observe that the result (i.e. if the condition in def. 2.4 is satisfied or not) is the same on almost all large random instances. The condition of def. 2.4 can hence be checked computationally with a small computer-time effort.

DEFINITION 2.5. Let  $\mathcal{N}_G$  be the number of satisfying assignments of an instance of the constraint satisfaction problem  $G$ . We define the annealed entropy  $s_{\text{ann}}$  to be

$$s_{\text{ann}} = \lim_{N \rightarrow \infty} \frac{1}{N} \log \mathbb{E}(\mathcal{N}_G), \quad (2.3)$$

where the expectation is over the graph ensemble. The quenched entropy is defined as

$$s_{\text{quen}} = \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}[\log(\mathcal{N}_G + 1)]. \quad (2.4)$$

DEFINITION 2.6. Let us also define the Bethe entropy [25, 27] that is associated to any BP fixed point as

$$s = \frac{1}{N} \sum_a \log(Z^a) - \frac{1}{N} \sum_i (l_i - 1) \log(Z^i). \quad (2.5)$$

where

$$Z^a = \sum_{\{s_i\}} \delta_{A_{\sum_i s_i, 1}} \prod_{i \in \partial a} \left( \prod_{b \in \partial i - a} \psi_{s_i}^{b \rightarrow i} \right), \quad (2.6a)$$

$$Z^i = \prod_{a \in \partial i} \psi_0^{a \rightarrow i} + \prod_{a \in \partial i} \psi_1^{a \rightarrow i}. \quad (2.6b)$$

The following statement stands on the basis of the cavity method: If a BP fixed point is asymptotically exact (for a given random graph ensemble) then the Bethe entropy (2.5) is equal to the quenched entropy (2.4).

The following result was obtained in [33] (section 5.1.2) for the occupation constraint satisfaction problems, and we conjecture it is general: The annealed entropy (2.3) is equal to the Bethe entropy evaluated in the uniform BP fixed point (when more than one uniform BP fixed point exists then consider the maximum of the Bethe entropy over the uniform fixed points).

If the problem is factorized and the uniform BP fixed point is asymptotically exact then the annealed entropy is equal to the quenched entropy,  $s_{\text{ann}} = s_{\text{quen}}$ . Our results in the remaining of this paper apply to random constraint satisfaction problems where indeed  $s_{\text{ann}} = s_{\text{quen}}$  (at least in some region of constraint densities). This condition is sometimes amenable to a rigorous proof, as it is in general weaker than  $\mathbb{E}(\mathcal{N}_G^2) < C [\mathbb{E}(\mathcal{N}_G)]^2$ , see [2]. If we are, however, interested in a fast heuristic check, then checking if the random constraint satisfaction problem is factorized may be more suitable.

**3. Basic properties of the random factorized locked problems.** For the locked problems, a detailed empirical analysis was done in [43, 42]. In this section we summarize the most relevant results (heuristically reasoned conjectures) of those works. It was found that a locked CSP is factorized in (at least) the two following cases:

- (a) **Any locked problem on random regular graphs**, that is when every variable is contained in  $L$  constraints. On regular graphs, the uniform fixed point of the BP equations then satisfies

$$\psi_0 = \frac{1}{Z} \sum_{r=0}^{K-1} \delta_{A(r), 1} \binom{K-1}{r} \psi_1^{(L-1)r} \psi_0^{(L-1)(K-1-r)}, \quad (3.1)$$

$$\psi_1 = \frac{1}{Z} \sum_{r=0}^{K-1} \delta_{A(r+1), 1} \binom{K-1}{r} \psi_1^{(L-1)r} \psi_0^{(L-1)(K-1-r)}, \quad (3.2)$$

where  $Z$  is the normalization. For the probability that a variable in occupied one has in this case in the  $N \rightarrow \infty$  limit

$$\chi = \frac{\psi_1^L}{\psi_1^L + \psi_0^L}, \quad (3.3)$$

Let us call  $x_r$  the probability that a constraint contains  $r$  occupied variables. Then

$$x_r = \frac{\binom{K}{r} \delta_{A(r),1} \psi_1^{r(L-1)} \psi_0^{(K-r)(L-1)}}{\sum_{t=0}^K \binom{K}{t} \delta_{A(t),1} \psi_1^{t(L-1)} \psi_0^{(K-t)(L-1)}}. \quad (3.4)$$

- (b) **The balanced locked problems** [42], are problems where the vector  $A$  is symmetric,  $A(i) = A(K-i)$  for all  $i = 0, \dots, K$  and this 0-1 symmetry is not spontaneously broken (that is when a satisfying assignment chosen uniformly at random has the same number of 0's and 1's up to a  $o(N)$  factor). Note that the absence of the symmetry breaking might depend on the degree distribution  $Q(l)$ . In the balanced locked problems the uniform BP fixed point  $\psi_1 = \psi_0 = \chi = 1/2$ . For the probability that a constraint contains  $r$  occupied variables we have here

$$x_r = \frac{\binom{K}{r} \delta_{A(r),1}}{\sum_{t=0}^K \binom{K}{t} \delta_{A(t),1}}. \quad (3.5)$$

A particularly simple case of (a) is the  $R$ -in- $K$  SAT where  $1 \leq R \leq K/2$ . If every variable has  $L$  connections and every constraint has to contain exactly  $R$  occupied variables, then the number of occupied variables is exactly  $MR/L$ , and thus  $\nu = R/K$ .

The authors of [43, 42] conjectured that when the  $N \rightarrow \infty$  limit of the Bethe entropy for a locked problem is positive, then the Bethe entropy is equal to the quenched entropy, and the BP fixed point reached from random initialization is asymptotically exact. If the Bethe entropy is negative then no satisfying assignment exists with high probability.

The Bethe entropy for all the balanced locked problems reads

$$s(\bar{l}) = \log 2 + \frac{\bar{l}}{K} \log \left[ 2^{-K} \sum_{r=0}^K \delta_{A(r),1} \binom{K}{r} \right], \quad (3.6)$$

where  $\bar{l}$  is the average degree of a variable (as we speak only about locked problems, the degree distribution has to have a zero weight on variables of degree zero and one). For all the locked problems on random regular (degree fixed to  $L$ ) graphs the entropy reads

$$s(L) = \frac{L}{K} \log \left[ \sum_{r=0}^K \delta_{A(r),1} \binom{K}{r} \psi_1^{(L-1)r} \psi_0^{(L-1)(K-r)} \right] - (L-1) \log [\psi_0^L + \psi_1^L], \quad (3.7)$$

where  $\psi_1, \psi_0$  is the fixed point of eqs. (3.1-3.2). This entropy simplifies further for the  $R$ -in- $K$  SAT on regular graphs (where the values of the  $\psi$ s obey the simple form discussed previously) where we get an explicit formula

$$s(L) = \frac{L}{K} \log \binom{K}{R} - (L-1) H \left( \frac{R}{K} \right), \quad (3.8)$$

where  $H(x) = -x \log x - (1-x) \log(1-x)$  is the entropy function.

The satisfiability transition  $l_s$  is then defined by

$$\text{satisfiability threshold } l_s : \quad s(l_s) = 0 \quad (3.9)$$

for the corresponding entropy function. For  $\bar{l} < l_s$  the problem has almost surely exponentially many solutions (the exponent being given by  $s(\bar{l})$ ) whereas for  $\bar{l} > l_s$  the problem almost surely does not have any solution.

The authors of [43, 42] also argued about the existence of a second phase transition in the locked problems,  $l_d < l_s$ , traditionally called in the physics literature the dynamical transition because of its connection to dynamics of glasses [32]. This critical point separates a region where for  $\{\sigma\}$  being a typical satisfying assignment the  $\psi_{\sigma_i}^{a \rightarrow i} = 1$ ,  $\psi_{-\sigma_i}^{a \rightarrow i} = 0$  defines a *stable* fixed point of the BP equations (2.1), from a region where this does not hold anymore. In other words, if an infinitesimal perturbation is introduced to these messages, the iteration of (2.1) goes back to the solution-related fixed point for  $l > l_d$ , but not for  $l < l_d$ . The authors of [43, 42] also conjectured that for  $\bar{l} > l_d$  a typical solution does not have solutions up to an extensive (i.e.  $\Theta(N)$ ) Hamming distance, whereas for  $\bar{l} < l_d$  there are other solutions at sub-extensive (i.e.  $o(N)$ ) Hamming distance. Let us call the phase corresponding to  $\bar{l} > l_d$  the separated phase, and the one corresponding to  $\bar{l} < l_d$  the non-separated phase.

For the locked problems on regular graphs, the following inequality always holds:  $2 < l_d < 3$ . In other words at  $L = 2$  the system is in the non-separated phase while for  $L \geq 3$  the solutions are always separated and the solution-corresponding fixed points are stable. For the balanced locked problems whenever the degree of every variable is larger or equal to three the system is in the phase where solutions are separated. When the fraction of variables of degree two is positive,  $Q(2) > 0$ , then the expression for  $l_d$  follows [42]:

$$\frac{l_d}{Q(2)} = 2(K-1) - 2 \frac{\sum_{r=1}^{K-2} r \binom{K-1}{r} \delta_{A(r+1),1} \delta_{A(r),0} \delta_{A(r-1),0}}{\sum_{r=0}^{K-2} \delta_{A(r+1),1} \binom{K-1}{r}}. \quad (3.10)$$

There is a deep connection between this dynamical threshold  $l_d$  and the reconstruction problem [10, 34, 26]. In the reconstruction problem one creates a tree with the same degree properties as the random graph. Then one considers a satisfying assignment chosen uniformly at random from all the possible ones. The reconstruction problem finally consists in deciding whether this assignment on the leaves of the tree contains some information about the value assigned to the root. In the locked problems the value of the root is always uniquely implied by the values of the leaves, as follows from the very definition of these problems. However, if an infinitesimal noise is introduced on the leaves then there is no information left if and only if  $\bar{l} < l_d$ . This value  $l_d$  was called the *small noise* reconstruction threshold in [42].

To summarize, the random locked factorized problems are in the non-separated phase for  $\bar{l} \leq l_d$ , which was shown to be algorithmically easy in [43, 42]. For  $l_d \leq \bar{l} \leq l_s$  the space of solutions is separated and it is then hard to find any solution. For  $\bar{l} \geq l_s$  no solution exists anymore.

**4. Equivalence of the random and planted ensembles.** The planted ensemble of graphs, which is the main subject of the present paper, is created in the following way:

TABLE 3.1

The critical values for all the balanced locked problems up to  $K = 8$  on the regular and truncated Poissonian ensembles. We remind here that the vector  $A$  codes for what are the allowed sums of variables around a constraint. We consider only problems where  $A(0) = A(K) = 0$ , that do not have a trivial all true or all false satisfying assignment. The integer value  $L_s$  (resp.  $L_l$ ) is defined as the first larger or equal to  $l_s$  (resp.  $l_l$ ), the stars denote that  $L_s = l_s$  (resp.  $L_l = l_l$ ). For definition of the threshold  $l_l$  see Sec. 6. The corresponding values of  $c$  are the coefficients that in the truncated Poisson distribution correspond to the average degree  $\bar{l}$ . The sign 'x' means that the problem ceases to be balanced before the instability arises.

A	$L_s$	$L_l$	$c_d$	$c_s$	$c_l$	$l_d$	$l_s$	$l_l$
00100	3	4*	1.256	1.853	2.821	2.513	2.827	3.434
0001000	4	6*	1.904	3.023	4.965	2.856	3.576	5.144
000010000	5	8*	2.337	3.942	6.994	3.116	4.276	7.039
5-in-10	5	10*	2.660	4.794	8.999	3.325	4.944	9.009
6-in-12	6	12*	2.918	5.455	11.00	3.502	5.586	11.00
01010	4*	$\infty$	1.904	3.594	$\infty$	2.856	4	$\infty$
0101010	6*	$\infty$	2.660	5.903	$\infty$	3.325	6	$\infty$
010101010	8*	$\infty$	3.132	7.978	$\infty$	3.654	8	$\infty$
0010100	6	46*	2.561	5.349	45.00	3.260	5.489	45.00
000101000	7	29*	2.975	6.650	28.00	3.542	6.708	28.00
001010100	8	> 100	3.110	7.797	> 100	3.638	7.822	> 100
010010010	6	x	2.173	4.896	x	3.014	5.083	x

TABLE 3.2

The critical values for all the regular (non-balanced) locked problems up to  $K = 6$ . The integer value  $L_s$  (resp.  $L_l$ ) is defined as the first larger or equal to  $l_s$  (resp.  $l_l$ ), the stars denote that  $L_s = l_s$  (resp.  $L_l = l_l$ ).

A	$L_s$	$L_l$
0100	3	3*
01000	3	4*
010000	3	5*
0100000	3	6*
001000	4	5*
0010000	4	6*

A	$L_s$	$L_l$
010100	5	> 50
0101000	6	> 50
010010	4	10
0100100	4	14
0100010	4	7

- (i) Make each of the  $N$  variables occupied with probability  $\chi$  (3.3), call the number of occupied variables  $N_1$ .
- (ii) Choose a degree sequence from the probability distribution  $Q(l)$  in such a way that  $KM = \sum_{i=1}^N l_i$ .
- (iii) For each constraint, and according to the probabilities  $x_r$  (3.4), choose the number  $r_a$  of occupied variables to which it is connected. Repeat until  $\sum_{a=1}^M r_a = \sum_{i=1}^{N_1} l_i$ . Here  $i$  are the indexes of the occupied variables. If this condition cannot be achieved go back to step (i) and repeat it until the condition is achievable.
- (iv) Now consider the  $r_a$  legs going out of every constraint  $a$ , order them arbitrarily and index them by  $i$  going from 1 to  $\sum_{a=1}^M r_a$ . Consider  $l_i$  legs going out from every occupied variable and index them. Choose a random permutation  $\pi$  of  $\sum_{a=1}^M r_a$  numbers, and connect the leg with index  $i$  going out from occupied variables to the leg with index  $\pi(i)$  going out from constraints. Do the same

with the empty variables and the remaining  $K - r_a$  legs going out from the constraints. Repeat until there are no double edges.

Note that there are several other models how to plant a solution (e.g. choose exactly the integer value of  $\chi N$  occupied variables in the step (i)), we could have chosen any other which is equivalent (for typical properties) to the above one in the  $N \rightarrow \infty$  limit.

**DEFINITION 4.1.** *Call a property of a large random graph drawn from a given random ensemble a thermodynamic property if and only if in the  $N \rightarrow \infty$  limit the probability that this property holds is smaller than  $1 - \exp(-cN)$ , where  $c$  is some  $\Theta(1)$  constant. In statistical physics of random systems it is often the case that large deviations are exponentially rare and hence all usually considered properties are thermodynamic in this sense. Without attempting a rigorous proof, in statistical physics the following examples are often assumed to be thermodynamic properties: The degree distribution, the entropy density, the fraction of occupied variables in a random satisfying assignment, the distance between two random satisfying assignments, etc. Properties that are not thermodynamic are all those relying on the behavior of exponentially rare instances, e.g. moments of some exponentially large quantities, as for instance the number of satisfying assignments.*

**DEFINITION 4.2.** *Consider two ensembles of random graphs  $A$  and  $B$ , we call the two ensembles asymptotically equivalent if and only if every property that is thermodynamic in ensemble  $A$  is thermodynamic in ensemble  $B$ , and vice versa.*

**DEFINITION 4.3.** *The planting is called quiet if the corresponding planted and random ensembles are asymptotically equivalent.*

Quiet planting intuitively means that if one is given a large random graph, one is not able to tell if that graph was drawn from the random or from the planted ensemble. This is because properties that are usually measured to distinguish the two graphs ensembles are thermodynamic (even if proving they are thermodynamic might be in general difficult).

The close relation between the random and the planted ensemble was explored in [2] for the random graph coloring and bi-coloring problems, see Theorem 6 and Theorem 7 in the appendix A of [2]. In statistical physics the quiet planting for graph coloring was discussed in [24].

**PROPOSITION 4.4.** *Denote  $\mathcal{N}_G$  the number of satisfying assignments of an instance  $G$  drawn from the random ensemble. If  $\mathbb{E}[\log(\mathcal{N}_G + 1)]/N$  is a thermodynamic property, i.e.*

$$\exists c > 0 : \forall \epsilon > 0 \lim_{N \rightarrow \infty} P(|\log \mathcal{N}_G - \mathbb{E}[\log(\mathcal{N}_G + 1)]| > \epsilon N) < e^{-cN} \quad (4.1)$$

and if the annealed entropy is equal to the quenched one, i.e.

$$\log \mathbb{E}(\mathcal{N}_G) = \mathbb{E}[\log(\mathcal{N}_G + 1)] + o(N). \quad (4.2)$$

then the planted ensemble and the random ensemble are asymptotically equivalent.

*Proof.* [of Prop. 4.4] In the random ensemble we are drawing graphs uniformly from all the graphs with a given degree distribution. In the planted ensemble we are drawing from the same set of graphs but with probability proportional to  $\mathcal{N}_G$ . Since (4.1) and (4.2) hold by assumption for the random ensemble, relation (4.1) holds for some  $c'$  also for the planted ensemble. Indeed, if (4.1) holds and if there would be larger than exponentially small probability that planting draws instances with  $|\log \mathcal{N}_G - \mathbb{E}[\log(\mathcal{N}_G + 1)]| > \epsilon N$  then (4.2) could not hold. For every other

thermodynamic property, i.e. such that large deviations are exponentially rare, the same argument applies.  $\square$

Statements equivalent to Prop. 4.4 first appeared in eq. (4) and Theorem 6 in [2]. Note that when the concentration condition (4.2) can be proven in a stronger form, then the condition on the exponentially rare large deviation can be weakened.

In statistical physics, using the cavity method arguments, the condition (4.2) can be evaluated. As we state at the end of sec. 2, when the Bethe entropy is asymptotically exact then the factorization and the equality of the quenched and annealed entropies are equivalent. Hence quiet planting is possible in all the factorized problems as long as the Bethe entropy is asymptotically exact. In the non-locked factorized problems, such as the random graph coloring, condition (4.2) ceases to be true strictly before the satisfiability threshold, as discussed in [24].

In the next section we argue that in the factorized locked problems (4.2) holds up to the satisfiability threshold, and hence the planted and the random ensembles are asymptotically equivalent for the factorized locked problems in the whole range of parameters corresponding to the satisfiable phase on the random ensemble.

**4.1. Second moment argument.** Relation (4.2) is in general rather hard to prove rigorously. Achlioptas and Coja-Oghlan [2] used instead a stronger condition  $\mathbb{E}(\mathcal{N}_G^2) < C [\mathbb{E}(\mathcal{N}_G)]^2$  which they proved for the coloring and the bi-coloring of factor-graphs problem for sufficiently sparse graphs.

For the factorized locked problems we conjecture that the relation  $\mathbb{E}(\mathcal{N}_G^2) < C [\mathbb{E}(\mathcal{N}_G)]^2$  holds in all the factorized locked problems on the purely random ensemble as long as  $\bar{l} \leq l_s$ .

The first and second moment of the number of solutions in the occupation problems has been computed for a general degree distribution in [42]. Based on numerical results it has been also argued non-rigorously in [42] that the above conjecture holds in the balanced locked problems. Here we illustrate that it also holds in the  $R$ -in- $K$  SAT on random  $L$ -regular graphs for  $L < l_s$ . The first moment entropy, defined by (2.3), is in the  $R$ -in- $K$  SAT on random  $L$ -regular graphs given by eq. (3.8). The second moment entropy  $s_{2\text{nd}} = \lim_{N \rightarrow \infty} \log \mathbb{E}(\mathcal{N}_G^2)/N$  is given by  $s_{2\text{dn}} = \max_t s_{2\text{nd}}(t)$  where [42]

$$s_{2\text{nd}}(t) = \frac{L}{K} \log \left\{ K! \sum_{s=0}^R \frac{\left[ \left(\frac{tR}{K}\right)^s \left[\frac{(1-t)R}{K}\right]^{2(R-s)} \left[1 + \frac{(t-2)R}{K}\right]^{K-2R+s} \right]^{1-\frac{1}{L}}}{(R-s)!(R-s)!s!(K-2R-s)!} \right\}. \quad (4.3)$$

The interpretation of the parameter  $0 \leq t \leq 1$  follows from expression

$$\mathbb{E}(\mathcal{N}_G^2) = \sum_{\sigma_1, \sigma_2} P(\sigma_1 \text{ SAT}, \sigma_2 \text{ SAT}), \quad (4.4)$$

where  $\sigma_1$  and  $\sigma_2$  are configurations and  $P(\cdot)$  is a probability over the graph ensemble. The parameter  $t$  in (4.3) is then the number of sites occupied in both  $\sigma_1$  and  $\sigma_2$  divided by number of sites occupied in one of the solutions,  $RN/K$ . We remind that in the  $R$ -in- $K$  SAT the satisfiability threshold is given by cancellation of the entropy (3.8)

$$l_s = \left[ 1 - \frac{\log\left(\frac{K}{R}\right)}{KH\left(\frac{R}{K}\right)} \right]^{-1}. \quad (4.5)$$

As  $s_{2\text{nd}}$  is a maximum of a function of a single variable  $t$ , we plot  $s_{2\text{nd}}(t)$  at  $L = l_s$  in Fig. 4.1. Evaluation of the polynomial function (4.3) for many values of  $R$  and  $K$  in Mathematica shows that for  $L < l_s$  we have  $2s_{\text{ann}} = s_{2\text{nd}} \geq 0$ , and for  $L > l_s$  we have  $s_{\text{ann}} = s_{2\text{nd}} \leq 0$ .

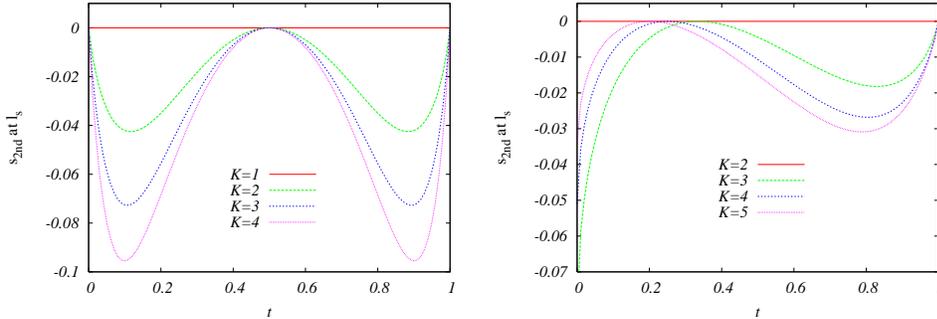


FIG. 4.1. The second moment entropy function  $s_{2\text{nd}}(t)$  (4.3) at  $l_s$  for several values of  $K$ , in the  $K$ -in- $2K$  SAT problem on the left, and 1-in- $K$  SAT on the right.

We also investigated numerically general formulas for the second moment presented in [42] and concluded that  $2s_{\text{ann}} = s_{2\text{nd}}$  for  $l < l_s$ , and  $s_{\text{ann}} = s_{2\text{nd}} \leq 0$  for  $l > l_s$ , holds also for all the other locked factorized problems.

**4.2. Satisfiable factorized locked problems equivalent to the planted ones.** We also conjecture that in the factorized locked models the planted ensemble is asymptotically equivalent to the ensemble of satisfiable instances in the whole region of  $\bar{l}$ .

For a general (non necessarily locked) constraint satisfaction problem the space of satisfying assignments is separated into clusters. We define the entropy  $s$  of a cluster as the logarithm of the number of assignments that belong to this cluster. We also define the complexity  $\Sigma(s)$  as the logarithm of the number of clusters of a given entropy  $s$ . The function  $\Sigma(s)$  is well defined even in the unsatisfiable region: when  $\Sigma(s) < 0$  it then corresponds to the large deviation function for the existence of a cluster of a given size [36]. It was argued in [24] that the cluster containing the planted configuration has a size  $s^*$  such that  $s^* = \text{argmax}[\Sigma(s) + s]$ . On the other hand from the large deviation interpretation of the  $\Sigma(s)$  function, most of the rare satisfiable instances in the unsatisfiable region will have one cluster of size  $s' = \text{argmax} \Sigma(s) \leq s^*$ .

In the locked problems, all clusters contain a single satisfying assignments, hence  $s = 0$  for all clusters. Keeping in mind the large deviation interpretation of the complexity  $\Sigma$  [36], the rare satisfiable instances have a single solution and should be asymptotically equivalent to planted instances. And hence the satisfiable and planted ensembles are asymptotically equivalent in the whole range of  $\bar{l}$  in the factorized locked problems.

**5. Single solution instances.** As discussed in the introduction, it is of practical importance to be able to create hard instances which have a single solution with a large probability. Based on the heuristic cavity method results of [24] we conjecture that in the region  $\bar{l} > l_s$  with high probability there is a single solution on large planted instances of the factorized locked problems (or a couple in case of balanced problems). In this section we prove (assuming the properties of the 1st and 2nd moment from the previous section) this statement for the  $R$ -in- $K$  SAT on random regular graphs.

We believe that the generalization of the proof is possible also for the other factorized locked problems.

First note that the first moment in the planted  $R$ -in- $K$  SAT  $s_{\text{ann,pl}} = \max_t s_{\text{ann,pl}}(t)$  is related in a simple way to the first and second moment in the purely random ensemble. It holds for the entropies

$$s_{\text{ann,pl}}(t) = s_{2\text{nd}}(t) - s_{\text{ann}}. \quad (5.1)$$

See an example of the function  $s_{\text{ann,pl}}(t)$  in Fig. 5.1.

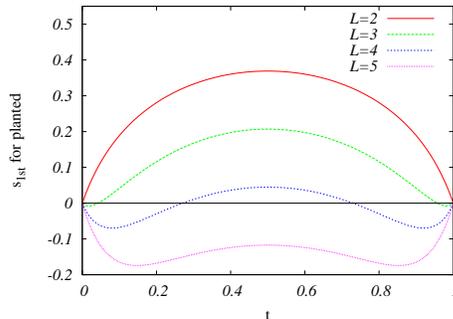


FIG. 5.1. The first moment entropy in the 4-in-8 SAT on  $L$  regular planted ensemble.

From the previous section it follows that for  $L > l_s$  the first moment entropy in the planted ensemble is a negative function for all  $0 < t < 1$ . The parameter  $t$  is in the planted ensemble interpreted as the distance from the planted solution. Therefore, for  $L > l_s$  there are no solutions at an extensive (i.e.  $\Theta(N)$ ) distance from the planted solution (except the solution at distance one in the balanced problems).

**THEOREM 5.1.** *Consider a large instance of the  $R$ -in- $K$  SAT problem drawn from the planted ensemble, the degree of variables be  $L > 2$ . Then there exists an  $\epsilon > 0$  such that with high probability there is no solution at distance smaller than  $\epsilon N$  from the planted solution. We will use the expander properties of regular bipartite graphs. The following theorem is well known in the theory of expanders [39].*

**THEOREM 5.2.** [Sipser and Spielman [39]] *Consider a random factor-graph with degree of variables  $L$  and degree of constraints  $K$ . Then, for any  $\delta < L-1$ , there exists a constant  $\epsilon > 0$ , such that with high probability for every set of  $\tilde{N} \leq \epsilon N$  variables the number of neighboring constraints is larger than  $\delta \tilde{N}$ . In other words the factor graph is a  $(\epsilon, \delta)$  expander.*

*Proof.* [of Theorem 5.1] Let us prove the statement by contradiction. Suppose that as  $N \rightarrow \infty$  for every  $\epsilon > 0$  there is a solution at distance smaller than  $\epsilon N$  from the planted solution. Denote the distance between the planted and this nearby solution  $N_1 = \epsilon' N$ . Now consider the factor-graph and the planted solution,  $N_1$  of variables have to be changed to reach the nearby solution. Since  $\epsilon'$  can be arbitrarily small Theorem 5.2 implies that there is at least  $\delta N_1$  constraints in which at least one variable has been changed. The property defining a locked constraint is that if a variable is changed then at least one other has to be changed in order to satisfy the constraint again. Hence each of the at least  $\delta N_1$  constraints have to be connected by at least two edges to the  $N_1$  changed variables. There is hence at least  $2\delta N_1$  edges connected to changed variables. The total degree of changed variables is  $LN_1$ , hence  $LN_1 > 2\delta N_1$ . But as  $\delta$  can be as near to  $L-1$  as we wish this inequality cannot

hold and we hence reached a contradiction. This proves that there exists  $\epsilon > 0$  such that with high probability there is no solution at distance smaller than  $\epsilon N$  from the planted one.  $\square$

Properties of the first moment in the planted ensemble together with Theorem 5.1 imply that in the planted  $R$ -in- $K$  SAT on random regular graphs there is almost surely a single solution (or a pair of solutions for  $R = K/2$ ).

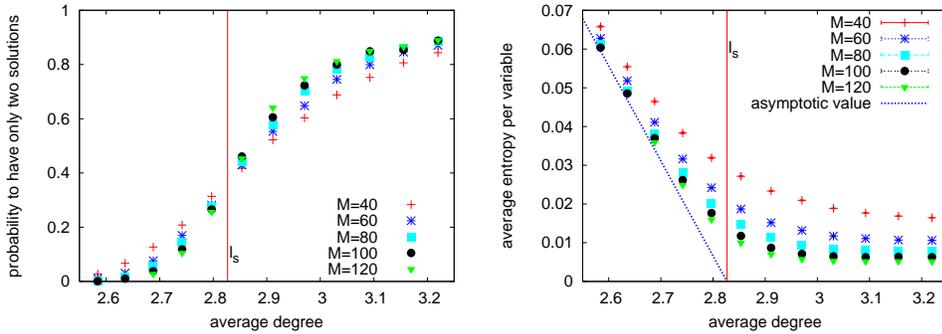


FIG. 5.2. Left: Probability (over 5000 instances) that there is a single pair of solutions in the 2-in-4 SAT as a function of the average degree and the size of the graph. Right: Data are the average entropy density (logarithm of the number of solutions per variables) of the instances. The line represents the entropy density in the  $N \rightarrow \infty$  limit, eq. (3.7). The data are obtained with the `relsat` algorithm [18]. In both parts we marked the threshold  $l_s = 2.827$ .

**6. Average computational hardness.** One of the most interesting aspects of the study of random constraint satisfaction problems is the average computational hardness of a given ensemble. This has been discussed extensively in both the computer science and the physics literature, in particular for the  $K$ -satisfiability and coloring problems. It has been shown empirically that the hardest instances lie very near to the satisfiability threshold  $l_s$ , and an easy-hard-easy pattern is often described [5, 31]. Later works focused on predicting up to which connectivity polynomial algorithms are able to find solutions, see e.g. [30, 41, 37]. Instances with a very large density of constraints are typically unsatisfiable. In some problems, e.g.  $K$ -satisfiability, no on average polynomial algorithms are known to show unsatisfiability for arbitrary large but constant density of constraints [6]. In other, more constraint, problems unit clause propagation based schemes were shown to be efficient [1]. In the planted instances, which are always satisfiable, it is known that for sufficiently large density of constraints solutions can be found in polynomial time, see e.g. [21, 8]. The situation in the planted factorized locked problems is very interesting: on top of the easy low and high constraint density phases we show that there also exists an intermediate hard phase, and we locate both the boundary thresholds.

We argued that in the satisfiable phase  $\bar{l} < l_s$  the planted and random ensembles are asymptotically equivalent, this includes the average behavior of algorithms. It was argued in [43, 42] that for average degree  $\bar{l} < l_d$  the locked problems are algorithmically easy whereas for  $l_d < \bar{l} < l_s$  they are on average hard.

The second hard-easy transition is particular to the planted ensemble and happens in the unsatisfiable phase  $\bar{l} > l_s$ . We will study the behavior of the BP equations initialized randomly to locate this transition.

**6.1. The spinodal point.** By definition of the factorized locked problems the belief propagation equations (2.1) initialized randomly converge to a uniform fixed point. But as the average degree is growing this ceases to be true. In the problems that we are studying here, there actually exists a critical average degree  $l_l$  beyond which belief propagation converges spontaneously towards the planted solution. This yields a clear hard-easy transition in the algorithmic complexity. In statistical physics terms this threshold  $l_l$  corresponds to a spinodal point of the liquid state [24]. The spinodal point also corresponds to the Kesten-Stigum bound [19, 20], and to the robust reconstruction threshold on trees [16]. This is yet another important connection between the reconstruction problem and the planted ensemble.

In order to compute the spinodal point let us first define matrix  $z(s'|s)$ . Consider a variable and one of its neighbors,  $z(s'|s)$  is then the probability that in the planted configuration the variable was assigned  $s'$  given that its neighbors was  $s$ . In the terms on reconstruction on trees  $z(s'|s)$  is the probability that in the broadcasting a variable was assigned  $s'$  given its parent was  $s$ . Components of  $z(s'|s)$  can be computed as

$$z(0|0) = \sum_{r=0}^K \left(1 - \frac{r}{K-1}\right) y_r(0), \quad z(1|0) = 1 - z(0|0), \quad (6.1)$$

$$z(1|1) = \sum_{r=0}^K \frac{r-1}{K-1} y_r(1), \quad z(0|1) = 1 - z(1|1). \quad (6.2)$$

where

$$y_r(0) = \frac{(K-r)x_r}{\sum_{t=0}^K (K-t)x_t}, \quad y_r(1) = \frac{r x_r}{\sum_{t=0}^K t x_t}, \quad (6.3)$$

where  $x_r$  is given by (3.4) or (3.5). Explicit formulas for the regular problems are

$$z(0|0) = \frac{\sum_{r=0}^{K-2} \binom{K-2}{r} \delta_{A(r),1} \psi_1^{r(L-1)} \psi_0^{(K-r-1)(L-1)}}{\sum_{r=0}^{K-1} \binom{K-1}{r} \delta_{A(r),1} \psi_1^{r(L-1)} \psi_0^{(K-r-1)(L-1)}}, \quad (6.4)$$

$$z(1|1) = \frac{\sum_{r=2}^K \binom{K-2}{r-2} \delta_{A(r),1} \psi_1^{(r-1)(L-1)} \psi_0^{(K-r)(L-1)}}{\sum_{r=1}^K \binom{K-1}{r-1} \delta_{A(r),1} \psi_1^{(r-1)(L-1)} \psi_0^{(K-r)(L-1)}}. \quad (6.5)$$

The first eigenvalue of this matrix is equal to one, and is associated with a trivial homogeneous eigenvector. The second eigenvalue of the matrix  $z$  is given by

$$\lambda = z(0|0) + z(1|1) - 1. \quad (6.6)$$

A well-known property of the reconstruction on a tree is that reconstruction is always possible beyond the so called Kesten-Stigum (KS) threshold [19, 20]. In our notation the KS condition says that if  $(L-1)(K-1)\lambda^2 > 1$  then the reconstruction is possible, i.e., the leaves asymptotically contain some information about the value sent by the root. In statistical physics the Kesten-Stigum condition is equivalent to the de Almeida-Thouless instability of the paramagnetic phase towards a spin-glass phase [9, 26, 22], that is for  $(L-1)(K-1)\lambda^2 > 1$  the belief propagation equations (2.1) do not converge. This can be seen from the fact that

$$\lambda = \frac{\partial \psi_1^{a \rightarrow i}}{\partial \psi_1^{b \rightarrow j}}, \quad (6.7)$$

where  $j \in \partial a \setminus i$ , and  $b \in \partial j \setminus a$ .

The eigenvalue  $\lambda$  and the condition for reconstructibility  $(L-1)(K-1)\lambda^2 > 1$  also appear in the problem of robust reconstruction on trees [16]. In the problem of robust reconstruction it is required that even if an arbitrary large fraction of the values on the leaves is erased there is still information about the root left.

The analysis of the instability of the uniform BP fixed point towards the planted solution then goes as follows. Consider a part of the factor-graph as depicted in Fig. 2.1. Denote the values of the messages in the uniform fixed BP fixed point by over-bars. Consider the incoming message to be perturbed from the uniform value as

$$\begin{pmatrix} \psi_1^{b \rightarrow j} = \overline{\psi}_1 + \epsilon \\ \psi_0^{b \rightarrow j} = \overline{\psi}_0 - \epsilon \end{pmatrix}. \quad (6.8)$$

Note that  $\epsilon$  can be both negative or positive. The equation (6.7) then implies that the outgoing message will be

$$\begin{pmatrix} \psi_1^{a \rightarrow i} = \overline{\psi}_1 + \lambda\epsilon \\ \psi_0^{a \rightarrow i} = \overline{\psi}_0 - \lambda\epsilon \end{pmatrix}. \quad (6.9)$$

In other words, any infinitesimal noise in one of the incoming message is multiplied by  $\lambda$  in the recursion.

We call the perturbation of the incoming message  $\epsilon_+$  if  $j$  was occupied in the planted configuration, and  $\epsilon_-$  otherwise. If the variable  $i$  was planted in the occupied state, then  $j$  was planted occupied with probability  $z(1|1)$ , and empty with probability  $z(0|1)$ . Similarly, if the variable  $i$  was planted in the empty state, then  $j$  was planted empty with probability  $z(0|0)$  and occupied with probability  $z(1|0)$ . Thus the evolution of the perturbation is governed by the equation:

$$\begin{pmatrix} \epsilon_+^{a \rightarrow i} \\ \epsilon_-^{a \rightarrow i} \end{pmatrix} = \lambda \begin{pmatrix} z(1|1) & z(0|1) \\ z(1|0) & z(0|0) \end{pmatrix} \begin{pmatrix} \epsilon_+^{b \rightarrow j} \\ \epsilon_-^{b \rightarrow j} \end{pmatrix}. \quad (6.10)$$

Moreover there are  $(K-1)(L-1)$  possible incoming messages in the regular graph, thus the criterion  $(K-1)(L-1)\lambda^2 = 1$ . If  $(K-1)(L-1)\lambda^2 < 1$  then the perturbation decreases and we find only the uniform BP fixed point, if on the contrary  $(K-1)(L-1)\lambda^2 > 1$  the uniform BP fixed point is unstable and a perturbation towards the planted configuration amplifies exponentially.

As the planted configuration corresponds to a stable BP fixed point<sup>1</sup> the BP iterations converge instead to the planted solution. Fig. 6.1 confirms that this is true even on rather small graphs. On the balanced locked problems, where we are not restricted to regular graphs, the correct condition is  $(K-1)\gamma\lambda^2 = 1$ , where  $\gamma$  is the mean of the excess degree distribution  $q(l) = (l+1)Q(l+1)/\bar{l}$ . The spinodal point  $l_l$ , see Tabs. 3.1,3.2, is then defined by

$$(K-1)(l_l-1)\lambda^2 = 1 \quad (6.11)$$

for the regular graphs, and

$$(K-1)\lambda^2 = \frac{1 - e^{-c_l}}{c_l} \quad (6.12)$$

---

<sup>1</sup>Note that in the above calculation we considered the stability around the uniform BP fixed point, if we consider the BP fixed point corresponding to the planted solution the perturbation does not amplify.

for the truncated Poissonian distribution.

The existence of this spinodal point, together with the conjecture about equivalence between the planted ensemble and the ensemble of satisfiable instances from the random ensemble, Sec. 4, implies that for  $\bar{l} > l_l$  it is easy to recognize almost all satisfiable instances of the locked problems. Similar conclusions, without a sharp threshold, were established for the coloring and satisfiability problems in [7, 11].

**6.2. Belief propagation as a solver.** Belief propagation reinforcement is a good solver in the region  $\bar{l} < l_d$  as shown empirically in [43, 42] in the random ensemble. Since the two ensembles are equivalent in that region, nothing changes for the planted ensemble. We have indeed verified this numerically.

Based on the above arguments, belief propagation equations converge to the uniform fixed point for  $\bar{l} < l_l$  and directly to the planted solution for  $\bar{l} > l_l$ . In order to verify that on finite size instances, we have performed the following numerical experiment: we have generated many planted instances for different sizes and average degrees (5000 instances for each set of parameters). We then iterated the BP equations (2.1) starting from random initial conditions. For numerical stability reasons we used dumping in the iterations, i.e. each time we computed a new message we kept one half of the sum of the new and old message. As a convergence criterion we used that the messages should not change more than  $2 \cdot 10^{-3}$  per message (we checked that a smaller criterion does not change the quality of results, and only slows down the computation). This way every iteration converged either to a configuration where the bias of each variable pointed towards the planted solution (or to its negation) or to a point very near to the uniform fixed point. Fig. 6.1 shows in what fraction of the runs we were able to find the planted solution and in particular it confirms that for  $\bar{l} > l_l$  it is easy to find it in linear time. On the right of the same figure we plot the average convergence time (given the criterion  $2 \cdot 10^{-3}$  per message). We see that around the spinodal point  $l_l$  the convergence time diverges from both the sides (slightly faster from the large degree side).

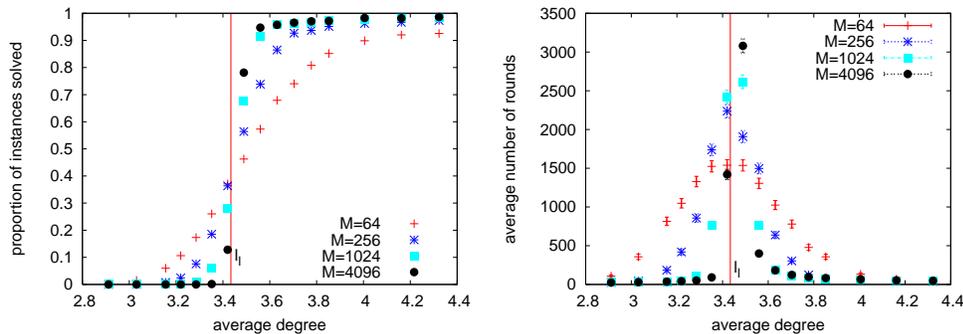


FIG. 6.1. *Belief propagation on the 2-in-4 SAT problem. Left: Probability that the belief propagation algorithm finds the planted configuration when initialized randomly plotted as a function of the average degree for several system sizes. Right: The convergence time dependence on the average degree. In both cases, we have stopped the BP iterations when the average change per message was less than  $2 \cdot 10^{-3}$ . In both parts we marked the spinodal threshold  $l_l = 3.434$ .*

**7. Conclusions and perspectives.** In this work we have studied a class of constraint satisfaction problems on a planted ensemble. The solution is planted in a *quiet* way, i.e. the planted configuration is one of the typical solutions of the resulting

instance. So far we know how to realize such plantings only on the factorized problems. We describe several connections between this quiet planting and the problem of reconstruction on trees.

We study the locked problems because of the simple structure on the space of their solutions — solutions are isolated points instead of clusters. This property makes the locked problems, however, very hard algorithmically. We focused on the class of occupation locked problems in this manuscript, all our results generalize easily to any factorized locked problem, on non-binary variables for example.

On non-locked but factorized problems, as e.g. graph coloring, the concept of quiet planting stays valid [24], however, the random and planted ensembles are not equivalent up to the satisfiability threshold. Moreover, in the unsatisfiable phase the planted ensemble has exponentially many solutions, instead of a single one as is the case in the locked problems. The non-locked problems are also much less friendly for first and second moment considerations. The phase diagram of the locked but non-factorized problem will not be very different from the one presented here. However, the thresholds will be different in the planted and random ensembles and the two ensembles are not equivalent.

One of the most important results of our work is the location of the algorithmically hard region, between  $l_d \leq \bar{l} \leq l_l$ , in the problems under investigation. It would be in particular interesting to design an algorithm which would provably find solutions in the region  $\bar{l} > l_l$ , as we have only heuristic and numerical arguments. This is also challenging in the non-locked problems, as e.g. graph coloring, where we predicted the planted Poisson ensemble to be easy above  $l_l = (q-1)^2$  (on planted regular graphs  $L_l = (q-1)^2 + 1$ ), where  $q$  is the number of colors. Results establishing that the planted ensemble on coloring is easy above  $Cq^2$ , where  $C$  is some constant quite larger than one, are already known [21, 8].

Finally, another consequence of our work worth discussing is that we know how to generate unique satisfying assignment instances - both in the hard and easy regions. Such instances are often used for evaluating the performance of the quantum annealing algorithm, but so far they have been generated with an exponential cost from an ensemble with unknown classical average computational complexity [40, 12]. In our opinion, these works should be repeated on instances of the locked problems. We conjecture that in the classically hard region also the quantum annealing will be exponential (this is because we anticipate a first order phase transition in the transverse magnetic field, as in [17]).

**Acknowledgments.** We thank to Dimitris Achlioptas, Amin Coja-Oghlan, Matti Jarvisalo, Andrea Montanari, Cris Moore and Guilhem Semerjian for fruitful discussions and suggestions.

#### REFERENCES

- [1] DIMITRIS ACHLIOPTAS, ARTHUR CHTCHERBA, GABRIEL ISTRATE, AND CRISTOPHER MOORE, *The phase transition in 1-in-k sat and nae 3-sat*, in *SODA '01: Proceedings of the twelfth annual ACM-SIAM symposium on Discrete algorithms*, Philadelphia, PA, USA, 2001, SIAM, pp. 721–722.
- [2] DIMITRIS ACHLIOPTAS AND AMIN COJA-OGHLAN, *Algorithmic barriers from phase transitions*. arXiv:0803.2122v2, 2008.
- [3] S. BEN-SHIMON AND D. VILENCHIK, *Message passing for the coloring problem: Gallager meets alon and kahale*, in *Proceedings of the 13th International Conference on Analysis of Algorithms*, DMTCS proc., 2007, pp. 217–226.
- [4] G. BIROLI AND M. MÉZARD, *Lattice glass models*, *Phys. Rev. Lett.*, 88 (2002), p. 025501.

- [5] PETER CHEESEMAN, BOB KANEFSKY, AND WILLIAM M. TAYLOR, *Where the Really Hard Problems Are*, in Proc. 12th IJCAI, San Mateo, CA, USA, 1991, Morgan Kaufmann, pp. 331–337.
- [6] VAŠEK CHVTAL AND ENDRE SZEMERDI, *Many hard examples for resolution*, J. ACM, 35 (1988), pp. 759–768.
- [7] A. COJA-OGHLAN, M. KRIVELEVICH, AND D. VILENCHIK, *Why almost all  $k$ -colorable graphs are easy to color*, in Proceedings of the 24th International Symposium on Theoretical Aspects of Computer Science (STACS), LNCS 4393, 2007, pp. 121–132.
- [8] AMIN COJA-OGHLAN, ELCHANAN MOSSEL, AND DAN VILENCHIK, *A spectral approach to analyzing belief propagation for 3-coloring*, Combinatorics, Probability and Computing, 18 (2009), pp. 881–912.
- [9] J. R. L. DE ALMEIDA AND D. J. THOULESS, *Stability of the Sherrington-Kirkpatrick solution of a spin-glass model*, J. Phys. A, 11 (1978), pp. 983–990.
- [10] WILLIAM EVANS, CLAIRE KENYON, YUVAL PERES, AND LEONARD J. SCHULMAN, *Broadcasting on trees and the Ising model*, Ann. Appl. Probab., 10 (2000), pp. 410–433.
- [11] R. MONASSON F. ALTARELLI AND F. ZAMPONI, *Can rare sat formulas be easily recognized? on the efficiency of message passing algorithms for  $k$ -sat at large clause-to-variable ratios*, J. Phys. A: Math. Theor., 40 (2007), pp. 867–886.
- [12] E. FARHI, J. GOLDSTONE, S. GUTMANN, J. LAPAN, A. LUNDGREN, AND D. PREDA, *A quantum adiabatic evolution algorithm applied to random instances of an np-complete problem*, Science, 292 (2001), p. 472.
- [13] U. FEIGE, E. MOSSEL, AND D. VILENCHIK, *Complete convergence of message passing algorithms for some satisfiability problems*, in Proc. of Random 2006, LNCS 4410, 2006, p. 339350.
- [14] ROBERT G. GALLAGER, *Low-density parity check codes*, IEEE Trans. Inform. Theory, 8 (1962), pp. 21–28.
- [15] R. G. GALLAGER, *Information theory and reliable communication*, John Wiley and Sons, New York, 1968.
- [16] SVANTE JANSON AND ELCHANAN MOSSEL, *Robust reconstruction on trees is determined by the second eigenvalue*, Ann. Probab., 32 (2004), pp. 2630–2649.
- [17] THOMAS JOERG, FLORENT KRZAKALA, JORGE KURCHAN, AND A. C. MAGGS, *Simple glass models and their quantum annealing*, Phys. Rev. Lett., 101 (2008), p. 147204.
- [18] R. J. BAYARDO JR. AND J. D. PEHOUSEK, *Counting models using connected components*, in Proc. 17th AAAI, Menlo Park, California, 2000, AAAI Press, pp. 157–162.
- [19] H. KESTEN AND B. P. STIGUM, *Additional limit theorems for indecomposable multidimensional galton-watson processes*, The Annals of Mathematical Statistics, 37 (1966), p. 1463.
- [20] ———, *Limit theorems for decomposable multi-dimensional galton-watson processes*, J. Math. Anal. Appl., 17 (1966), p. 309.
- [21] M. KRIVELEVICH AND D. VILENCHIK, *Semi-random models as benchmarks for coloring algorithms*, in Proceedings of the Third Workshop on Analytic Algorithmics and Combinatorics (ANALCO), 2006, pp. 211–221.
- [22] FLORENT KRZAKALA, ANDREA MONTANARI, FEDERICO RICCI-TERSENGHI, GUILHEM SEMERJIAN, AND LENKA ZDEBOROVÁ, *Gibbs states and the set of solutions of random constraint satisfaction problems*, Proc. Natl. Acad. Sci. U.S.A, 104 (2007), p. 10318.
- [23] F. KRZAKALA, M. TARZIA, AND L. ZDEBOROVÁ, *A Lattice Model for Colloidal Gels and Glasses*, Phys. Rev. Lett., 101 (2008), p. 165702.
- [24] FLORENT KRZAKALA AND LENKA ZDEBOROVÁ, *Hiding quiet solutions in random constraint satisfaction problems*, Phys. Rev. Lett., 102 (2009), p. 238701.
- [25] F. R. KSCHISCHANG, B. FREY, AND H.-A. LOELIGER, *Factor graphs and the sum-product algorithm*, IEEE Trans. Inform. Theory, 47 (2001), pp. 498–519.
- [26] MARC MÉZARD AND ANDREA MONTANARI, *Reconstruction on trees and spin glass transition*, J. Stat. Phys., 124 (2006), pp. 1317–1350.
- [27] M. MÉZARD AND A. MONTANARI, *Information, Physics, Computation*, Oxford University Press, Oxford, 2009.
- [28] M. MÉZARD AND G. PARISI, *The bethe lattice spin glass revisited*, Eur. Phys. J. B, 20 (2001), p. 217.
- [29] M. MÉZARD, G. PARISI, AND M. A. VIRASORO, *Spin-Glass Theory and Beyond*, vol. 9 of Lecture Notes in Physics, World Scientific, Singapore, 1987.
- [30] M. MÉZARD, G. PARISI, AND R. ZECCHINA, *Analytic and algorithmic solution of random satisfiability problems*, Science, 297 (2002), pp. 812–815.
- [31] DAVID G. MITCHELL, BART SELMAN, AND HECTOR J. LEVESQUE, *Hard and easy distributions for SAT problems*, in Proc. 10th AAAI, Menlo Park, California, 1992, AAAI Press, pp. 459–465.

- [32] A. MONTANARI AND G. SEMERJIAN, *On the dynamics of the glass transition on bethe lattices*, J. Stat. Phys., 124 (2006), pp. 103–189.
- [33] T. MORA, *Géométrie et inférence dans l'optimisation et en théorie de l'information*, PhD thesis, Université Paris-Sud, 2007. <http://tel.archives-ouvertes.fr/tel-00175221/en/>.
- [34] ELCHANAN MOSSEL, *Reconstruction on trees: Beating the second eigenvalue*, Ann. Appl. Probab., 11 (2001), pp. 285–300.
- [35] J. PEARL, *Reverend bayes on inference engines: A distributed hierarchical approach*, in Proceedings American Association of Artificial Intelligence National Conference on AI, Pittsburgh, PA, USA, 1982, pp. 133–136.
- [36] O. RIVOIRE, *The cavity method for large deviations*, J. Stat. Mech., (2005), p. P07004.
- [37] SAKARI SEITZ, MIKKO ALAVA, AND PEKKA ORPONEN, *Focused local search for random 3-satisfiability*, J. Stat. Mech., (2005), p. P06006.
- [38] BART SELMAN, HENRY A. KAUTZ, AND BRAM COHEN, *Local search strategies for satisfiability testing*, in Proceedings of the Second DIMACS Challenge on Cliques, Coloring, and Satisfiability, Michael Trick and David Stifler Johnson, eds., Providence RI, 1996.
- [39] MICHAEL SIPSER AND DANIEL A. SPIELMAN, *Expander codes*, IEEE Transactions on Information Theory, 42 (1996), pp. 1710–1722.
- [40] A. P. YOUNG, S. KNYSH., AND V.N. SMELYANSKIY, *Size dependence of the minimum excitation gap in the quantum adiabatic algorithm*, Phys. Rev. Lett., 101 (2008), p. 170503.
- [41] L. ZDEBOROVÁ AND F. KRZAKALA, *Phase transitions in the coloring of random graphs*, Phys. Rev. E, 76 (2007), p. 031131.
- [42] L. ZDEBOROVÁ AND M. MÉZARD, *Constraint satisfaction problems with isolated solutions are hard*, J. Stat. Mech., (2008), p. P12004.
- [43] ———, *Locked constraint satisfaction problems*, Phys. Rev. Lett., 101 (2008), p. 078702.