

Phase transitions in number theory: From the birthday problem to Sidon setsBartolo Luque,¹ Iván G. Torre,¹ and Lucas Lacasa^{2,*}¹*Departamento de Matemática Aplicada y Estadística, ETSI Aeronáuticos, Universidad Politécnica de Madrid, Spain*²*School of Mathematical Sciences, Queen Mary University of London, Mile End Road, London E1 4NS, United Kingdom*

(Received 8 August 2013; revised manuscript received 25 September 2013; published 12 November 2013)

In this work, we show how number theoretical problems can be fruitfully approached with the tools of statistical physics. We focus on g -Sidon sets, which describe sequences of integers whose pairwise sums are different, and propose a random decision problem which addresses the probability of a random set of k integers to be g -Sidon. First, we provide numerical evidence showing that there is a crossover between satisfiable and unsatisfiable phases which converts to an abrupt phase transition in a properly defined thermodynamic limit. Initially assuming independence, we then develop a mean-field theory for the g -Sidon decision problem. We further improve the mean-field theory, which is only qualitatively correct, by incorporating deviations from independence, yielding results in good quantitative agreement with the numerics for both finite systems and in the thermodynamic limit. Connections between the generalized birthday problem in probability theory, the number theory of Sidon sets and the properties of q -Potts models in condensed matter physics are briefly discussed.

DOI: [10.1103/PhysRevE.88.052119](https://doi.org/10.1103/PhysRevE.88.052119)

PACS number(s): 05.70.Fh, 02.10.De, 89.20.-a, 89.75.-k

I. INTRODUCTION

From the celebrated work on the critical behavior of random Boolean satisfiability (random k -SAT) [1,2], we have seen how statistical physics and computer science have built bridges between each other, as it has been recognized that the theory underlying optimization problems [3] and the physics of disordered systems [4] shares many similarities [5,6]. In computer science, one can investigate the situations under which decision problems are satisfiable, and study the geometry and accessibility to those solutions. Decision problems can indeed be interpreted under a statistical physics formalism, where the cost function to be minimized (in random k -SAT, this is for instance the number of violated constraints) relates to the Hamiltonian of a disordered system at zero temperature: in this latter situation, the physical system tries to adopt the ground state or minimal energy configuration, only reachable in some circumstances, while in others frustration (unsatisfiability) can develop. In the last years, an exciting multidisciplinary environment has witnessed the efforts of describing optimization problems within the physics of disordered problems, including replica-symmetry-breaking solutions in combinatorial problems [4] or the description of phase transitions (threshold phenomena) in decision problems [3,7–9]. In a nutshell, the mutual interchange of approaches and techniques from physics to computer science and viceversa have proved to be a valuable input in both fields [6].

Here we explore the possibility of further extending this fruitful relation to number theory. Can arithmetics and number theory be seen as a natural system, subject to scientific scrutiny much in the form physics observes physical reality? Can numbers be considered as units that interact locally according to some arithmetic properties and hence be amenable to a statistical mechanics description? Several works following this epistemological approach range from the reinterpretation of the nontrivial Riemann ζ zeros as the eigenspectrum of a quantum Hamiltonian [10] to the onset of phase transitions

in number partitioning problems [5,9]. Within complexity science, a fundamental question is to find the minimal amount of ingredients a system needs to possess to evidence emergent behavior. Number theoretic systems can indeed be thought of as the utterly purest, unadorned and simplest models where complexity may develop [11], and therefore constitute a privileged playground for scientific research. In this paper we show how number theoretical problems are susceptible to be approached from a computer science and statistical physics perspective, and indeed show nontrivial emergent properties. For concreteness, we focus on the arithmetics of Sidon sets [12] and show that this number theoretical statement is susceptible to be treated as a random Constraint Satisfaction Problem (rCSP), and show links with both the statistical physics of disordered systems and with some classical problems in probability theory. In doing so, we will unveil the onset of a phase transition in the satisfiability of such system—a so-called zero-one law in the mathematical jargon [13,14]—which we will show is analytically tractable.

II. RANDOM g -SIDON: A NUMBER-THEORETIC DECISION PROBLEM

In number theory, a set of k different positive integers $\mathcal{X} = \{S_1, S_2, \dots, S_k\}$, is a so-called Sidon set [12] if all the sums of two elements $S_i + S_j$ from the set (where $i, j = 1, \dots, k$) are different (except when they coincide because of commutativity). For example, $\{1, 2, 5, 10, 16, 23\}$ is a Sidon set, whereas $\{1, 3, 7, 10, 17, 23\}$ is not Sidon since $1 + 23 = 7 + 17$. Sidon sets recurrently appear in different areas of mathematics including Fourier analysis, group theory, or number theory. An extension of Sidon sets allows in the definition for g repetitions, accordingly, \mathcal{X} is called g -Sidon provided that any sum of two elements $S_i + S_j$ is *repeated* at most $g - 1$ times (note that when $g = 1$, a g -Sidon set reduces to a Sidon set). From its first beginnings [15], number theorists were interested in the extremal properties of Sidon sets, concretely in calculating upper and lower bounds of the maximal size of g -Sidon sets formed from an integer interval $[1, M]$ for diverging values of M , a topic and focus which still has an

*lucas.lacasa2@gmail.com

intense research activity [16–19]. Less attention (if any) has been paid in the onset of zero-one laws in such systems [20]. Here we show that a statistical physics approach is helpful in this task. We begin by recasting the concept of Sidon sets and propose a rCSP approach as it follows: let $[1, M]$ be a pool of positive integers and extract at random from this pool a set of k different numbers $\mathcal{X} = \{S_1, S_2, \dots, S_k\} \subset [1, M]$. Which is the (satisfaction) probability $P_g(k, M)$ that \mathcal{X} is a g -Sidon set? We will call this rCSP the random g -Sidon decision problem. First, notice that this problem can be rephrased in a statistical physics formalism in the following terms: given the initial set \mathcal{X} , proceed to build a secondary set $\tilde{\mathcal{X}} = \{S_i + S_j\}_{i,j=1,\dots,k} = \{\tilde{S}_n\}_{n=1,\dots,k(k+1)/2} \subset [2, 2M]$ formed by all the sums of two elements from the set \mathcal{X} . As the sums $S_i + S_j$, $i = 1, \dots, k$ are allowed, $\tilde{\mathcal{X}}$ will be formed by $k(k+1)/2$ positive integers between 2 and $2M$. $\tilde{\mathcal{X}}$ is therefore the output of a Sidon set \mathcal{X} if every number in \tilde{S}_n is distinct. Let us define the number of matches of \tilde{S}_n in $\tilde{\mathcal{X}}$ as

$$h_n = \sum_{m=1, m \neq n}^{\frac{k(k+1)}{2}} \delta[\tilde{S}_n - \tilde{S}_m], \quad (1)$$

where the Kronecker function $\delta[x] = 1$ if $x = 0$ and 0 otherwise. Then, for $g = 1$, we can define the following function \mathcal{H} :

$$\mathcal{H}(g = 1, k, M, \tilde{\mathcal{X}}) = \sum_{n=1}^{\frac{k(k+1)}{2}} h_n. \quad (2)$$

\mathcal{H} computes “the degree of Sidonlikeness” of a configuration $\{S_j\}$, i.e., the total number of matches in $\{\tilde{S}_n\}$. Indeed, $\mathcal{H} = 0$ for Sidon sets while $\mathcal{H} > 0$ for non-Sidon sets. Stated as a rCSP, $\mathcal{H} = 0$ and $\mathcal{H} > 0$ distinguish the satisfiable and unsatisfiable phases, respectively, such that satisfiability is reached for configurations that minimize this function. \mathcal{H} can be seen as the physical internal energy of the statistical mechanics system $\tilde{\mathcal{X}}$, where each configuration of the variables $\{\tilde{S}_n\}$ is a given microstate with energy \mathcal{H} . Each of the $k(k+1)/2$ “spins” can take discrete values in $[2, 2M]$. If random fluctuations of the spin values were allowed, the equilibrium properties of this system at temperature T would be given by the Boltzmann measure in $[2, 2M]$ $\mu(\tilde{\mathcal{S}}) = \exp[-\beta\mathcal{H}(\tilde{\mathcal{S}})]/Z$, where $\beta \sim 1/T$, and Z is the normalization (partition) function. In general, the system will be in the phase (Sidon/non-Sidon) that minimizes the Helmholtz free energy $F = \mathcal{H} - TS$. Since no random fluctuations of the values of the variables are allowed in our system, the system is to be considered at zero temperature, and variables will try to occupy the ground state energy, that is, the configuration that minimizes \mathcal{H} , being this the satisfiable phase if $\min[\mathcal{H}] = 0$, and getting frustrated in the unsatisfiable one if the minimum energy state available is larger than zero.

The case of g -Sidon sets ($g > 1$) is slightly more involved:

$$\mathcal{H}(g, k, M, \tilde{\mathcal{X}}) = \sum_{n=1}^{\frac{k(k+1)}{2}} (h_n - g + 1)\theta[h_n - g], \quad (3)$$

where the Heaviside step function $\theta[x] = 0$ if $x < 0$ and $\theta[x] = 1$ if $x \geq 0$. Hence, g -Sidon sets and non- g -Sidon comply $\mathcal{H} = 0$ and $\mathcal{H} > 0$, respectively.

III. NUMERICAL RESULTS: FROM CROSSOVER TO PHASE TRANSITIONS

The order parameter which naturally associates with sat/unsat phases in the random g -Sidon problem is the satisfaction probability $P_g(k, M)$, that describes the probability of a randomly extracted set of k elements from $[1, M]$ to be g -Sidon, i.e., to fulfill $\min[\mathcal{H}] = 0$. We will consider k as the control parameter and in what follows we firstly explore numerically the behavior of $P_g(k, M)$ as a function of M and g . The behavior for $g = 1$, as a result of (ensemble-averaged) Monte Carlo simulations, is shown in Fig. 1. First, notice that the transition between satisfiability [$P_g(k, M) \approx 1$] and unsatisfiability [$P_g(k, M) \approx 0$] occurs at increasing values $k(M, g)$. In order the control parameter to be intensive, we rescale it as $\alpha = k/k_c(M, g)$, where we make use the standard in percolation theory and set $P_g(k_c, M) = 1/2$. In the bottom inset panel of the same figure we show the explicit dependence $k_c(M)$, which we find agrees with the expression $k_c(M) \sim M^{1/4}$. Notice that k is not actually extensive (linear) in M and therefore the ratio k/M typically used in the satisfiability theory [3] does not work here.

In the upper inset panel of the same figure we also show the behavior $P_1(\alpha, M)$ as a function of α , for different pool cardinals M . The collapse of the order parameter under a single smooth curve points out that the behavior is independent of M , which means that the pool’s size does not play any relevant role, and $P_1(k, M) \xrightarrow{k \rightarrow \alpha} P_1(\alpha)$. Also, such transition is smooth both for finite sizes and in the thermodynamic limit ($k \rightarrow \infty$, $M \rightarrow \infty$, α finite), that is, the system seems to evidence a simple crossover and no threshold phenomenon occurs.

In order to cast light in the effect of g , in Fig. 2 we show the numerical results of $P_g(k, M)$ for a concrete pool size

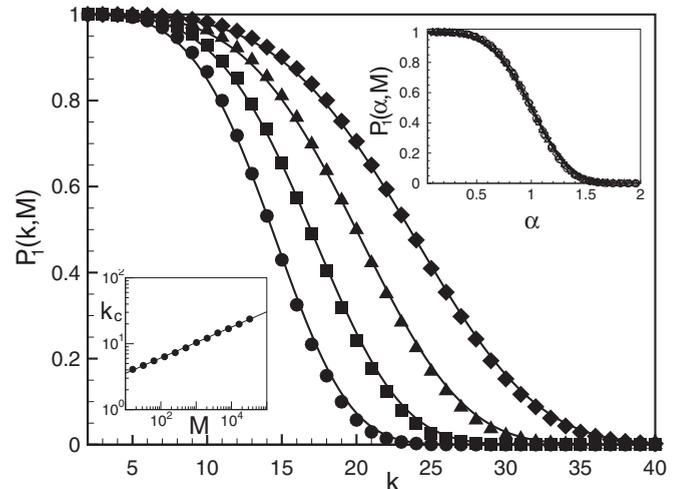


FIG. 1. (Central panel) Numerical simulations of satisfaction probability $P_1(k, M)$ as a function of k , for $M = 2^{12}$ (circles), $M = 2^{13}$ (squares), $M = 2^{14}$ (triangles), $M = 2^{15}$ (diamonds), averaged over 10^4 realizations. Solid lines are the prediction of Eq. (5). (Inset bottom panel) Log-log plot of the transition point $k_c(M)$, defined as $P_1(k_c, M) = 1/2$, showing the scaling $k_c \sim M^{1/4}$. (Inset top panel) Collapsed satisfaction probability $P_1(\alpha, M)$, for $\alpha = k/k_c$, finding a continuous sigmoidal function independent of M . Solid line is a prediction of the theory $P_1(\alpha) = 2^{-\alpha^4}$ (see the text).

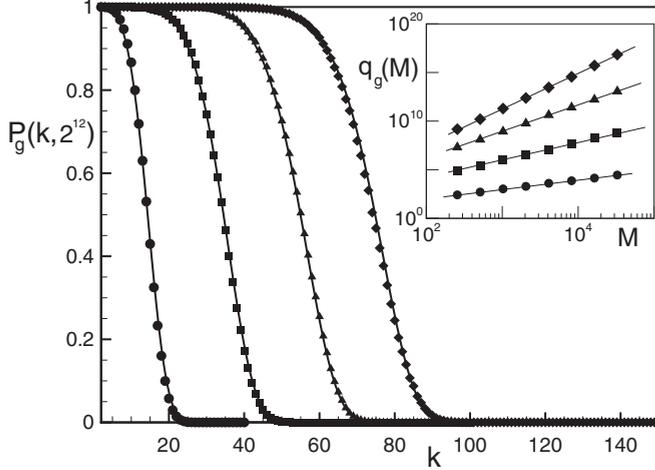


FIG. 2. (Central panel) Numerical simulations of satisfaction probability $P_g(k, M = 2^{12})$ as a function of k , for $g = 1$ (circles), $g = 2$ (squares), $g = 3$ (triangles), and $g = 4$ (diamonds) (Monte Carlo simulations averaged over 5×10^4 realizations). Solid lines are the prediction of the theory (see the text). (Inset panel) Numerical scalings of $q_g(M)$ for different values of g , finding in each case a power law relation albeit with different slopes, suggesting the two-variable scaling $q_g(M) = A(g)M^{r(g)}$. The expression q_g is related to k_c through the transformation $q_g(M) \equiv k_c^{2(g+1)} / ((g+1)! 2^{2g+1} \log 2)$ (see the text). The specific shapes of $A(g)$ and $r(g)$ are plotted the insets of Fig. 3.

$M = 2^{12}$, for different values of g . Again, the transition point k_c shows a dependency not only with M but also with g (see the inset panel). When the control parameter is properly made intensive, and at odds with the phenomenology for $g = 1$, we find that the transition *sharpens* for increasing values of g . This result is further confirmed in Fig. 3, where we plot the behavior of $P_g(\alpha, M)$ for different values of M and g . For finite g the satisfaction probability adopts again a universal M -independent sigmoidal curve $P_g(k, M) \xrightarrow{k \mapsto \alpha} P_g(\alpha)$, although this curve gets sharper around $\alpha = 1$ as g increases. This sharpening further suggests that g plays the role of an *effective* system size, such that the crossover that takes place for finite g seems to develop into a true phase transition in the limit of $g \rightarrow \infty$. This is a genuinely counterintuitive result having in mind that in decision problems the apparent system's size is usually related to pool's size, here M , which in our case is an irrelevant variable.

IV. ANALYTICAL DEVELOPMENTS

In what follows we support this phenomenology with some analytical calculations. Incidentally, note at this point that if \tilde{S}_i were drawn uniformly from $[2, 2M]$, the 1-Sidon problem would be equivalent to the celebrated birthday problem [21], a standard in probability theory that calculates the probability that, in a set of k randomly chosen people, not a single pair will share the same birthday (with a year containing $N = 2M - 1$ days). In that case Eq. (2) would also be equivalent to the Hamiltonian of a N -Potts model in a mean field approximation (no explicit space) widely used in solid state physics [22]. Similarly, if again \tilde{S}_i were drawn uniformly, the g -Sidon problem would reduce in the secondary description to the

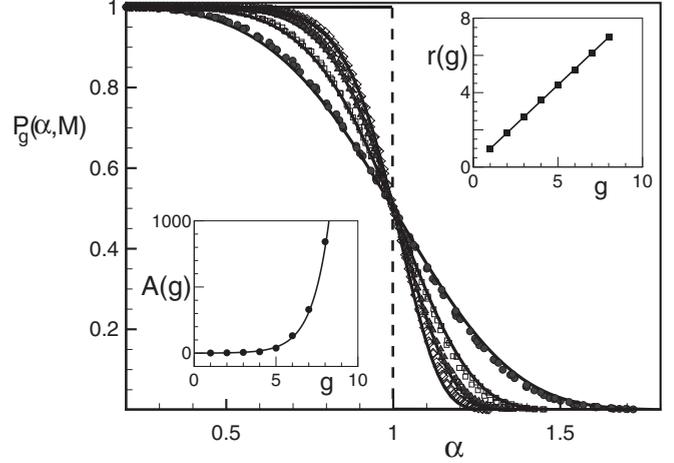


FIG. 3. Rescaled satisfaction probability $P_g(\alpha, M)$ as a function of α , for several values of $g = 1$ (gray circles), $g = 2$ (squares), $g = 3$ (gray triangles), and $g = 4$ (diamonds), and $M = 2^{12}, 2^{13}, 2^{14}, 2^{15}$. For each g , there is a universal collapse curve $P_g(\alpha)$ which turns to be independent of M , however, such universal curves get sharper around $\alpha = 1$ for increasing values of g , suggesting the onset of a phase transition in the limit $g \rightarrow \infty, \alpha$ finite. Solid lines for finite g and the Heaviside function for $g \rightarrow \infty$ are predictions of the theory (see the text). (Inset top panel) Linear fit of the correcting exponent $r(g) \approx 0.85g + 0.13$ ($R^2 = 0.9998$). (Bottom inset panel) Exponential fit of the correcting factor $A(g) \approx 0.4 \exp(0.95g)$ ($R^2 = 0.9985$).

so-called generalized birthday problem [23], which calculates the probability that if k people are selected at random, $g + 1$ people will not have the same birthday (or alternatively, the probability that at most g people will share the same birthday). In our problem \tilde{S}_n are the birthdays and $2M - 1$ the days in one year. Since \tilde{S}_n are nonuniformly sampled in $[2, 2M]$ (as they are the result of the sum pairs in $\{S_i\}$), we will only take the birthday problem/Potts model as naive approximations of the random g -Sidon problem.

Suppose that a year has N days quote $P_g(n, N)$ the probability that no $g + 1$ people, of n people selected at random, have the same birthday. Then the following recursive relation holds approximately:

$$P_g(n + 1, N | n, N)^N \approx P_{g-1}(n, N).$$

Using Bayes' theorem and taking logarithms in the preceding equation, we find

$$\log P_g(n + 1, N) - \log P_g(n, N) = \frac{1}{N} \log P_{g-1}(n, N).$$

Now, for sufficiently large $N \gg n$, the left-hand side in the latter expression is a first-order discretization of $\partial \log P_g / \partial n$. In the continuum limit this yields a partial differential equation for the evolution of P_g

$$\frac{\partial P_g(n, N)}{\partial n} = \frac{P_g(n, N)}{N} \log P_{g-1}(n, N)$$

that along with initial condition $P_g(0, N) = 1$ for all g has the following solution:

$$P_g(n, N) = \exp\left(-\frac{n^{g+1}}{(g+1)!N^g}\right). \quad (4)$$

The generalized birthday problem has an analytically unmanageable closed-form expression [23], that nonetheless has been treated asymptotically by some authors [24,25]. We find noticeable that the asymptotic solutions to the problem, derived using slightly sophisticated combinatorial and statistics techniques, agrees with our expression in Eq. (4), obtained following a straightforward argument.

If we assume in our problem that the values \tilde{S}_n are uncorrelated and result of independent trials in $[2, 2M]$, $P_g(k, M)$ results from the change of variables: $N \rightarrow 2M$, $n \rightarrow k(k+1)/2 \approx k^2/2$:

$$P_g(k, M) = \exp\left(-\frac{k^{2(g+1)}}{(g+1)!2^{2g+1}A(g)M^{r(g)}}\right), \quad (5)$$

where we have formally substituted M^g with $A(g)M^{r(g)}$ because the $n = k(k+1)/2$ variables are correlated in our system. In order to quantitatively compare our theory with finite-size numerics, we shall express higher order deviations from this equation introducing a correcting exponent $r(g) \neq g$, whose first-order perturbative expansion reads $r(g) = r_0 + r_1g$, and similarly for the normalizing factor $A(g)$. The concrete values of the free parameters are then found using a simple self-consistent argument, imposing that the correct scalings $k_c(M, g)$ shall be found at $P_g(k_c, M) = 1/2$. After a little algebra we find that $q_g(M) \equiv k_c^{2(g+1)}/((g+1)!2^{2g+1} \log 2) = A(g)M^{r(g)}$, where the specific fits for $A(g)$ and $r(g)$ are shown in the inset panels of Fig. 3. This expression reduces to $k_c \sim M^{1/4}$ for $g = 1$, on agreement with previous numerical evidence (inset panel of Fig. 1). The predicted values of $P_g(k, M)$ are accordingly plotted in solid lines along with the numerics in Figs. 1–3 for different values of M and g , showing an excellent agreement in every case.

Incidentally, our theory also predicts that the maximal size k_{\max} of a g -Sidon set—a classical question in number theory—should satisfy $\binom{M}{k_{\max}}P_g(k_{\max}, M) = 1$, whose leading order for $g = 1$ is $k_{\max} = \mathcal{O}(M^{\sqrt{2}-1})$, on agreement with a recent theorem by Ruzsa [16] (see Fig. 4).

Finally, as a function of the intensive control parameter $\alpha = k/k_c$, Eq. (5) reduces to

$$P_g(\alpha, M) = 2^{-\alpha^{2(g+1)}}. \quad (6)$$

It is important to highlight that this law holds independently of the concrete values of $A(g)$ and $r(g)$ —that is, it is a direct consequence of the theory—and also holds without needs to impose taking the limit $M \rightarrow \infty$ (see Fig. 3). The solution that we obtain is a universal sigmoid function for finite g [in the case of $g = 1$, the curve is $P_1(\alpha) = 2^{-\alpha^4}$, on excellent agreement

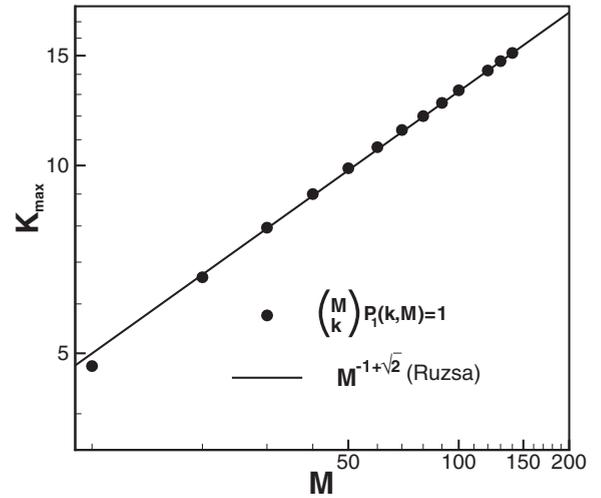


FIG. 4. (Dots) Maximal size of a 1-Sidon set as a function of M , as predicted by the phenomenological theory underlying the random Sidon decision problem. (Solid line) Maximal 1-Sidon set bound given by Ruzsa [16].

with the numerics in the inset panel of Fig. 1]. This sigmoid sharpens for increasing values of g , and in the thermodynamic limit $g \rightarrow \infty$ (what necessarily implies $k \rightarrow \infty$), α finite a zero-one law emerges

$$P_\infty(\alpha, M) = \begin{cases} 1 & \text{if } \alpha < 1 \\ 0 & \text{if } \alpha > 1 \end{cases}. \quad (7)$$

That is, whereas the decision problem only evidences a crossover for all finite g , this transition indeed becomes abrupt and converts to a true phase transition in the thermodynamic limit.

V. CONCLUSION

To conclude, we have shown how the methods and focus of statistical physics and theoretical computer science can be fruitfully applied in the realm of number theory. We have found and described, both numerically and analytically, a previously unnoticed phase transition within the properties of g -Sidon sets, with the exotic peculiarity that M —the analog of the number of possible values of each spin, for instance q in the q -Potts model—does not play any relevant role in the onset of the phase transition, while the finite-size role is played here by the combinatorial parameter g . The extension of these approaches to other number theoretical problems, and the establishment of new links between these fields are important open problems to be further addressed.

[1] P. Cheeseman, B. Kanefsky, and W. M. Taylor, in *Where the Really Hard Problems Are*, Proceedings of IJCAI-91, edited by J. Mylopoulos and R. Rediter (Morgan Kaufmann, San Mateo, CA, 1991), p. 331.
 [2] S. Kirkpatrick and B. Selman, *Science* **264**, 1297 (1994).
 [3] *Handbook of Satisfiability*, edited by A. Biere, M. Heule, H. van Maaren, and T. Walsh (IOS Press, Fairfax, VA, 2009).

[4] M. Mezard, G. Parisi, and M. Virasoro, *Spin Glass Theory and Beyond* (World Scientific, Singapore, 1987).
 [5] S. Mertens, *Theoret. Comput. Sci.* **265**, 79 (2001).
 [6] C. Moore and S. Mertens, *The Nature of Computation* (Oxford University Press, Oxford, 2011).
 [7] R. Monasson, R. Zecchina, S. Kirkpatrick, B. Selman, and L. Troyansky, *Nature (London)* **400**, 133 (1999).

- [8] R. Monasson and R. Zecchina, *Phys. Rev. Lett.* **76**, 3881 (1996).
- [9] S. Mertens, *Phys. Rev. Lett.* **81**, 4281 (1998).
- [10] M. V. Berry and J. P. Keating, *SIAM Rev.* **41**(2), 236 (1999).
- [11] B. Luque, O. Miramontes, and L. Lacasa, *Phys. Rev. Lett.* **101**, 158702 (2008).
- [12] K. J. Compton, in *0-1 Laws in Logic and Combinatorics, Algorithms and Order*, edited by I. Rival, NATO ASI Series, (Kluwer Academic Publishers, Dordrecht, 1988), pp. 353–383.
- [13] S. Shelah and J. Spencer, *J. Amer. Math. Soc.* **1**, 97 (1988).
- [14] H. L. Abbott, *Canad. Math. Bull.* **33**, 335 (1990).
- [15] P. Erdos and P. Turan, *J. London Math. Soc.* **16**, 212 (1941).
- [16] I. Z. Ruzsa, *J. Number Theory* **68**, 63 (1998).
- [17] J. Cilleruelo, I. Z. Ruzsa, and C. Trujillo, *J. Number Theory* **97**, 26 (2002).
- [18] G. Yu, *J. Number Theory* **122**, 211 (2007).
- [19] K. O’Bryant, *Electron. J. Combin.* **11**, 1 (2004).
- [20] T. Luczak and J. Spencer, *J. Amer. Math. Soc.* **4**, 451 (1991).
- [21] W. Feller, *An Introduction to Probability Theory and Its Applications*, 3rd ed., Vol. 1 (John Wiley & Sons, New York, 1970).
- [22] F. Y. Yu, *Rev. Mod. Phys.* **54**, 235 (1982).
- [23] E. H. Mckinney, *Am. Math. Mon.* **73**, 385 (1966).
- [24] H. Mendelson, *J. Comb. Theory A* **30**, 351 (1981).
- [25] N. Henze, *Statistics & Probability Letters* **39**, 333 (1998).