



Article

# The Self-Similarity of Industrial SAT Instances

Carlos Ansótegui<sup>1</sup>, Maria Luisa Bonet<sup>2</sup>, Jesús Giráldez-Cru<sup>3</sup> and Jordi Levy<sup>4,\*</sup><sup>1</sup> DIEI, Universitat de Lleida (UdL), Spain; carlos@diei.udl.cat<sup>2</sup> LSI, Universitat Politècnica de Catalunya (UPC), Spain; bonet@cs.upc.edu<sup>3</sup> DaSCI, DECSAI, Universidad de Granada (UGR), Spain; jgiralde@ugr.es<sup>4</sup> Artificial Intelligence Research Institute (IIIA-CSIC), Spain; levy@iiia.csic.es

\* Correspondence: levy@iiia.csic.es; Tel.: +34 935809570

Version October 26, 2020 submitted to Mathematics

**Abstract:** In the last years, we have witnessed a remarkable progress of algorithms solving Boolean satisfiability (SAT). The success of these algorithms has been especially relevant in a large number of industrial or real-world applications, for which these SAT solvers are nowadays an essential core part of their solving processes. Interestingly enough, these applications include a very diverse and heterogeneous range of domains, such as hardware verification, planning, and cryptography, among others. Unfortunately, the reasons of the good performance of these solvers on this variety of industrial benchmarks are not completely understood yet. Since SAT solvers' efficiency is fundamental in various domains, obtaining a better understanding of these algorithms and the reasons of their good performance is crucial.

In order to shed light on this question, SAT solvers are often viewed as *complex systems* with many interconnected components (e.g., conflict analysis and learning mechanism, database management, search restarts) interacting in many unpredictable ways. There is the common belief that the resulting emergent behavior of these complex systems takes advantage of a certain underlying structure of the SAT formula, which is shared by the majority of these industrial problems regardless the domain they come from. Recently, there have been some attempts of characterizing this structure under the lens of complex networks, with the purpose of better understanding the success of the solvers, and potentially improving them.

In this paper, we analyze the structure of industrial SAT instances under the lens of self-similarity, and study how the execution of SAT solvers affect that structure. Many real-world graphs exhibit *self-similar structure* (with small *fractal dimension*), which means that after *rescaling* (replacing groups of nodes by a single node), the same kind of structure can be observed. In our analysis, in which we represent SAT instances as graphs, we observe that many industrial SAT formulas exhibit the same kind of structure. Moreover, we analyze how this structure evolves by effects of learning new clauses during the search. In particular, we observe that learned clauses usually contain variables that are *close* in the graph representation of the formula. This is, the learning mechanism tends to work *locally*. On the contrary, this learning mechanism on random SAT formulas –which do not exhibit any structure at all– is unable to generate these local clauses. This difference contributes to explain the success of modern SAT solvers on industrial problems.

**Keywords:** satisfiability; SAT; fractal dimension; complex networks.

## 1. Introduction

The Boolean Satisfiability Problem (SAT) is the problem of deciding whether the Boolean variables of a propositional formula can be assigned in such a way that the formula is evaluated as true. A SAT formula is said to be satisfiable if such an assignment does exist, and unsatisfiable otherwise. This problem is central in Computer Science, from both theory and applications perspectives. It is the first NP-complete problem, which means that existing algorithms can run during exponentially long executions, in the worst case. But some algorithms –the so known *modern SAT solvers*– show a very good performance solving a certain kind of instances:



instances encoding problems from real-world applications or industrial domains, such as hardware and software verification, scheduling, cryptography, or planning, among others [1–3].

It is well established that those industrial SAT instances have a very distinct nature to randomly generated SAT formulas. As a consequence, SAT solving algorithms more *specialized* in random benchmarks do not usually perform well on industrial ones, whilst solvers for industrial purposes usually show a very bad performance on random SAT formulas. For instance, these *industrial-specialized* solvers are able to solve very large industrial SAT instances with millions of variables in a few seconds, whereas they are unable to solve random SAT formulas with hundreds of variables in a few hours or even days.

Modern SAT solvers implement the Conflict-Driven Clause Learning (CDCL) algorithm [4]. Built on the basis of the classical DPLL algorithm [5,6], which is a depth-first search algorithm exploring the variables assignment space of the formula, CDCL incorporates a number of complex techniques such as conflict analysis and clause learning [7], sophisticated variable branching heuristics [8], clause removal policies to delete useless learned clauses [9], random restarts of the search [10], or pre-/in-processing techniques to detect variables dependencies and simplify the formula [11].

Since CDCL SAT solvers can be seen as complex systems where there are many unpredictable interactions between its components, the reasons that explain why these solvers show such a different performance on random and industrial instances, or more specifically, the reasons that explain their success on a large variety of industrial benchmarks, are not completely understood yet. It is therefore fundamental to shed light to this open question, which is the main motivation of this manuscript. A better understanding of these solvers and the reasons of their success will possibly help to improve them.

The main component of CDCL SAT solvers is clause learning [12,13]. These learned clauses are the result of the conflicts found during the search, and they prevent the solver to find the same conflicts in the future. Moreover, a very extended heuristic implemented in many CDCL SAT solvers is VSIDS [8]. This heuristic selects the next branching variable based on the already found conflicts. This results into a synergy between expanding the search (branching variables) and delimiting it (found conflicts); this synergy is somehow *captured* by learned clauses.<sup>1</sup> Therefore, learned clauses play a key role in the algorithm.<sup>2</sup>

A common intuition to explain the distinct performance of CDCL SAT solvers between random and industrial benchmarks comes from the belief that these solvers are able to exploit a certain underlying structure of the formula [9,16–20], that is shared by the majority of the benchmarks. In our work, we are inspired by the work of *complex networks* where the general structure of real-world graphs is studied. An extended observation is that the classical Erdős-Rényi random graph model [21] is not suitable to study the structure of these complex networks. In order to characterize the underlying structure of industrial SAT instances, we follow a similar approach. In particular, we represent SAT instances as graphs, and then analyze some graph features in these graph representations. For this purpose, we use two graph models. In one model, SAT formulas are represented as bipartite graphs where their nodes represent either the variables or the clauses of the formulas, and an edge between a variable-node and a clause-node indicates the existence of such a variable in such a clause. In the second model, nodes only represent variables, and there is an edge between two variable-nodes if there is a clause containing both variables.

Following this approach of representing SAT instances as graphs, an interesting observation is described in [22]. In this work, it is studied the community structure of industrial SAT formulas, and authors observe that most of the industrial formulas show a very clear community structure. In a graph with clear community structure, it can be found a partition of disjoint sets of nodes, called *communities*, such that the graph has more edges connecting nodes of the same community. In the context of SAT instances, this means that the variables of the formula mostly appear in clauses with the same subset of variables. Interestingly enough, they observe that

---

<sup>1</sup> Every variable has an *activity* score, equally initialized. When a variable occurs in a conflict, it increases its activity, and it is very likely that such a variable occurs in the resulting learned clause. The VSIDS heuristic selects the variable with highest activity. Since this variable probably occurs in existing learned clauses, it is very likely it triggers new conflicts. This is the synergy between the conflicts and the heuristic, *captured* by learned clauses.

<sup>2</sup> It is intriguing if they also play the same role in other common heuristics, as LRB [14] or VMTF [15].

78 the quality of this community structure tends to decrease by the addition of learned clauses.<sup>3</sup> It remains unclear,  
79 however, whether this decrease is achieved *quickly* or *slowly*.

80 In this work, we focus on the analysis of the *self-similarity* of SAT formulas and their *fractal dimension*.  
81 In a graph with self-similar structure, the structure of the graph is preserved at different *scales* (for instance,  
82 by grouping sets of nodes into a single nodes). It also means that the diameter  $d^{max}$  of the graph grows as  
83  $d^{max} \sim n^{1/d}$ , being  $d$  its fractal dimension. Notice that in random graphs or small-world graphs  $d^{max} \sim \log n$ .

84 Our experimental evaluation shows that many application SAT instances exhibit self-similar structure with  
85 a small fractal dimension ranging between 2 and 4. We conjecture that this structure, as well as other common  
86 properties in these instances, such as community structure [22] or scale-free structure [23], are already existent  
87 in many of the high-level problems encoded into the SAT instances. Thus, for example, hardware-verification  
88 problems may exhibit this structure as a consequence of the circuits they encode. We also evaluate how the  
89 fractal dimension is affected by the effect of learning new clauses during the search. In particular, we observe  
90 that most learned clauses contain variables that are already related by other clauses. In the graph representation  
91 of the formula, this means that learned clauses create edges between already connected nodes (i.e., variables) or  
92 very *close* ones. In other words, the learning mechanism tends to work *locally*. This observation allows us to  
93 determine that the underlying structure of the formula is destroyed *slowly*.

94 This paper is an extended and revised version of [24]. Some preliminary results were also included in Jesús  
95 Giráldez-Cru's PhD thesis [25].

96 The rest of the paper proceeds as follows. Some preliminary concepts are introduced in Section 2. We study  
97 the self-similar structure of graphs and SAT formulas in Sections 3 and 4 respectively. We analyze how this  
98 structure is affected by the CDCL SAT solving techniques in Section 5. Finally, related works and conclusions  
99 are respectively described in Sections 6 and 7.

## 100 2. Preliminaries

101 Given a set of Boolean variables  $X = \{x_1, \dots, x_n\}$ , a *literal* is an expression of the form  $x_i$  or  $\neg x_i$ . A  
102 *clause*  $c$  of size  $s$  is a disjunction of  $s$  literals,  $l_1 \vee \dots \vee l_s$ . We note  $s = |c|$ , and say that  $x \in c$ , if  $c$  contains the  
103 literal  $x$  or  $\neg x$ . A *CNF formula* or *SAT instance* of length  $t$  is a conjunction of  $t$  clauses,  $c_1 \wedge \dots \wedge c_t$ . A *k-CNF*  
104 *formula* is a conjunction of  $k$ -sized clauses.

105 An undirected graph is a pair  $(V, E)$  where  $V$  is a set of vertices and  $E \subseteq V \times V$  is a set of edges between  
106 the vertices. The degree of a vertex  $x$  is defined as  $\deg(x) = |\{y \in V | (x, y) \in E\}|$ . A bipartite graph is a tuple  
107  $(V_1, V_2, E)$  where  $E \subseteq V_1 \times V_2$ .

108 We represent SAT instances using two graph models. In the Variable Incidence Graph model (VIG, for  
109 short), vertices represent variables, and edges represent the existence of a clause relating two variables. A clause  
110  $x_1 \vee \dots \vee x_n$  results into  $\binom{n}{2}$  edges, one for every pair of variables.

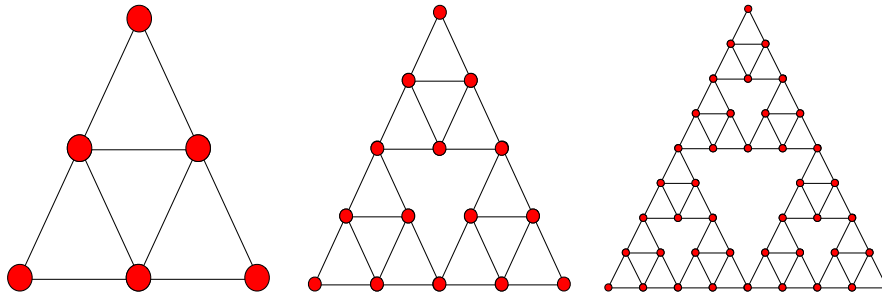
**Definition 1** (Variable Incidence Graph (VIG)). *Given a SAT instance  $\Gamma$  over the set of variables  $X$ , its variable incidence graph is a graph  $(X, E)$  with set of vertices the set of Boolean variables, and the set of edges  $E$ :*

$$(x, y) \in E \Leftrightarrow c \in \Gamma \wedge x \in c \wedge y \in c \quad (1)$$

111 In the Clause-Variable Incidence Graph model (CVIG, for short), vertices represent either variables or  
112 clauses, and edges represent the occurrence of a variable in a clause.

---

<sup>3</sup> The addition of learned clauses tends to destroy the community structure of the formula.



**Figure 1.** Example of a graph with self-similar structure at distinct scales.

**Definition 2** (Clause-Variable Incidence Graph (CVIG)). *Given a SAT instance  $\Gamma$  over the set of variables  $X$ , its clause-variable incidence graph is a bipartite graph  $(X, \{c \mid c \in \Gamma\}, E)$ , with vertices the set of variables and the set of clauses, and a set of edges  $E$ :*

$$(x, y) \in E \Leftrightarrow y \in \Gamma \wedge x \in y \quad (2)$$

113 These models do not fully represent the SAT formula, ignoring some of its features. For instance, having  
 114 the CVIG it cannot be known the sign of a variable in a clause, or the length of clauses cannot be known having  
 115 the VIG. In the literature there are other graph models for SAT formulas capturing all their information. For  
 116 instance, the factor graph (similar to the CVIG) with two kind of edges (one for positive literals and the other for  
 117 negative ones). However, these simplified models have already all the information required for our analysis.

### 118 3. The Self-Similar Structure of Graphs

119 In this section, we describe the concepts and techniques used to analyze the self-similar structure of graphs.  
 120 In particular, we introduce the Burning by Node Degree (BND) algorithm (see Alg. 1), which will be used to  
 121 analyze the fractal dimension in very large graphs, as the graph representation of industrial SAT formulas are.  
 122 We also analyze the relation between fractal dimension and diameter.

123 In Figure 1 we depict a graph with clear self-similar structure in order to graphically describe the concept of  
 124 fractal dimension of a graph.

125 In order to define the fractal dimension of a graph, we follow the principle of self-similarity. We first need  
 126 to define the concepts of distance between nodes and graph coverage. For this purpose, we will use the definition  
 127 of box covering [26].

128 **Definition 3.** *The **distance** between two nodes of a graph is the minimum number of its edges we need to follow  
 129 to go from one node to the other.*

130 **Definition 4.** *The **diameter**  $d^{max}$  of a graph is the maximal distance between any two nodes of the graph.*

131 **Definition 5.** *Given a graph  $G$ , a **box**  $B$  of size  $l$  is a subset of its nodes such that the distance between any pair  
 132 of them is strictly smaller than  $l$ .*

133 **Definition 6.** *The **minimum number of boxes** of size  $l$ ,  $M(l)$ , to cover the graph is the minimum number of  
 134 boxes such that every node of the graph belongs to some box.*

135 Notice that  $M(1)$  is equal to the number of nodes of  $G$ , and  $M(d^{max} + 1)$  is the number of connected  
 136 components of the graph (i.e.,  $M(d^{max} + 1) = 1$  if the graph is fully connected). Notice also that the function  
 137  $M(l)$  has the same value for any  $l$  strictly greater than  $d^{max}$  (i.e., the number of connected components of the  
 138 graph).

139 **Definition 7.** *A graph has the **self-similarity** property if the function  $M(l)$  decreases polynomially, i.e.  $M(l) \sim$   
 140  $l^{-d}$ , for some value  $d$ . In this case, we call  $d$  the **dimension** of the graph.*

141 In graphs with self-similar structure, this function  $M(l)$  decreases *slowly*, in contrast to other graphs for  
 142 which this functions decreases faster. Therefore, in this latter case of graphs, the diameter is comparatively  
 143 smaller than in those graphs showing fractal dimension.

144 **Lemma 8.** *Computing the function  $M(l)$  is NP-hard.* <sup>4</sup>

145 **Proof.** We prove that computing  $M(2)$  is already NP-hard by reducing the graph coloring problem to the  
 146 computation of  $M(2)$ . Given a graph  $G$ , let  $\bar{G}$ , the complement of  $G$ , be a graph with the same nodes, and where  
 147 any pair of distinct nodes are connected in  $\bar{G}$  iff they are not connected in  $G$ . Boxes of size 2 in  $\bar{G}$  are cliques,  
 148 thus they are sets of nodes of  $G$  without an edge between them. Therefore, the minimal number of colors needed  
 149 to color  $\bar{G}$  is equal to the minimal number of cliques needed to cover  $\bar{G}$ , i.e.  $M(2)$ .  $\square$

150 There are several efficient algorithms that approximate  $M(l)$ . They compute upper bounds of  $M(l)$ . They  
 151 are called *burning* algorithms (see [27]). Following a greedy strategy, at every step they try to select the box that  
 152 covers (burns) the maximal number of uncovered (unburned) nodes. Although they are polynomial algorithms,  
 153 we still need to do some further approximations to make the algorithms of practical use in very large graphs.

154 First, instead of boxes, we will use *circles*.

155 **Definition 9.** *Given a graph  $G$ , a **circle** of radius  $r$  and center  $c$  is a subset of its nodes such that the distance  
 156 between any of them and the node  $c$  is strictly smaller than  $r$ .*

157 **Definition 10.** *The **minimum number of circles** of radius  $r$ ,  $N(r)$ , to cover the graph is the minimum number  
 158 of circles such that every node of the graph belongs to some circle.*

159 Notice that any circle of radius  $r$  is inside of a box of size  $2r - 1$  (the opposite is in general false) and any  
 160 box of size  $l$  is inside a circle of radius  $l$  (it does not matter what node of the box we use as center). Notice also  
 161 that every radius  $r$  and center  $c$  characterizes a unique circle.

162 According to Hausdorff's dimension definition,  $N(r) \sim r^{-d}$  also characterizes self-similar graphs of  
 163 dimension  $d$ . We can approximate this fractal dimension using the *Maximum-Excluded-Mass-Burning (MEMB)*  
 164 algorithm [27], which works as follows: Consider a graph  $G$  and a radius  $r$ . We compute an upper bound of the  
 165 number of circles with radius  $r$  necessary to cover the graph  $N(r)$ . We start with all nodes set to unburned. At  
 166 every step, for every possible node  $c$ , we compute the number of unburned nodes covered by the circle of center  
 167  $c$  and radius  $r$ , then select the node  $c$  that maximizes this number, and burn the new covered nodes. Starting with  
 168  $r = 1$ , we repeat this process for every value of  $r$  until we get  $N(r) = 1$ .

169 The MEMB algorithm is still too costly for our purposes. Notice that in every step a node is selected as a  
 170 center, the algorithm has previously visited every node of the graph in order to find the best center. So in order to  
 171 make the algorithm more efficient, we apply the following algorithm, called *Burning by Node Degree (BND)*,  
 172 and described in Alg. 1. First, we compute  $N(1)$  as the number of nodes in the graph (line 1). Then, we iterate to  
 173 compute  $N(r)$  with  $r > 1$ , until we reach the number of connected components of the graph (line 3). In each  
 174 iteration, we use node degrees to order the set of nodes:  $\langle c_1, \dots, c_n \rangle$  such that  $degree(c_i) \geq degree(c_j)$ , when  
 175  $i < j$  (line 7). Now, for  $i = 1$  to  $n$ , if  $c_i$  has not been burned in previous iterations, the circle of center  $c_i$  and  
 176 radius  $r$  is selected, burning all the nodes contained in this circle (line 8-10). This is repeated until all nodes are  
 177 burned, obtaining  $N(r)$  for this value of  $r$  (line 11). Notice that this criterion may not maximize the number of  
 178 burned nodes.

179 We recall that BND (see Alg. 1) computes upper bounds of the function  $N(r)$  until it gets the number of  
 180 connected components of the graph. Notice that  $N(d^{max} + 1) = 1$ , placing the center of this circle in any node

---

<sup>4</sup> In [27] it is stated the same result, but they prove the wrong reduction. They reduce the computation of  $M(2)$  to the graph coloring problem.

**Algorithm 1:** Burning by Node Degree (BND)

---

```

Input: Graph  $G = (V, E)$ 
Output: vector[int]  $N$ 
1  $N[1] := |V|;$ 
2  $\text{int } i := 2;$ 
3 while  $N[i - 1] > \text{connectedComponents}(G)$  do
4   vector[bool]  $\text{burned}(|V|, \text{false});$  // all nodes initially unburned;
5    $N[i] := 0;$ 
6   while  $\text{existsUnburnedNode}(\text{burned})$  do
7      $c := \text{highestDegreeUnburnedNode}(G, \text{burned});$ 
8      $S := \text{circle}(c, i);$  // circle with center  $c$  and radius  $i$ ;
9     foreach  $x \in S$  do
10     $\text{burned}[x] := \text{true};$ 
11     $N[i] := N[i] + 1;$ 
12   $i := i + 1;$ 

```

---

181 of the path between the two most distant nodes. However, it is possible that BND computes the approximation of  
 182  $N(r)$  for values of  $r$  greater than  $d^{\max} + 1$ .

183 In the next section, the accuracy and performance of algorithms MEMB and BND will be compared. This  
 184 analysis will allow us to justify the use of the latter in our experimental evaluation.

185 In this paper we analyze the function  $N(r)$  for the graphs obtained from a SAT instance following the VIG  
 186 and CVIG models. These two functions are denoted  $N^V(r)$  and  $N^C(r)$ , respectively, and they relate to each  
 187 other as follows.

188 **Lemma 11.** *If  $N^V(r) \sim r^{-d}$  then  $N^C(r) \sim r^{-d}$ .  
 If  $N^V(r) \sim e^{-\beta r}$  then  $N^C(r) \sim e^{-\frac{\beta}{2} r}$ .*

189 **Proof.** Notice that, for any formula, given a circle of radius  $r$  in the VIG model, using the same center and radius  
 190  $2r - 1$  we can cover the same variable nodes in the CVIG model. With radius  $2r$  we can also cover all clauses  
 191 adjacent to some covered variable. Hence  $N^C(2r) \leq N^V(r)$ .

192 Conversely, given a circle of radius  $2r$  in the CVIG model, we consider two possibilities. If the center is a  
 193 variable node, we cover the same variables in the VIG model using a circle of radius  $r$  and the same center. If the  
 194 center is a clause  $c$ , to cover the same variables in the VIG model, we need a circle of radius  $r + 1$  centered in a  
 195 variable node adjacent to  $c$ . Hence  $N(r + 1) \leq N^C(2r)$ .

196 Therefore  $N^V(r + 1) \leq N^C(2r) \leq N^V(r)$ , and  $N^V(r) \sim N^C(2r)$ . From this asymptotic relation, we  
 197 can derive the two implications stated in the lemma.  $\square$

198 From the previous lemma, it can be concluded that if a SAT instances has a (perfect) self-similar structure,  
 199 its fractal dimension is the same in both models VIG and CVIG. In this case, we would obtain a line if we  
 200 represent  $N(r)$  as a function of  $r$  in double-logarithmic scale, being  $-d$  its slope. If there is an exponential decay  
 201 in the function  $N(r)$ , then the decay factor in the VIG is double of the decay factor in the CVIG model. In that  
 202 case, representing  $N(r)$  in semi-logarithmic axes, we would obtain a line with slope  $-\beta$ . In our analysis, we  
 203 always represent  $N(r)$  in double-logarithmic axes. Therefore, if  $N(r)$  decays exponentially, we will observe a  
 204 concave curve.

### 205 3.1. Fractal Dimension versus Diameter

206 The function  $N(r)$  determines the *maximal radius*  $r^{\max}$  of a connected graph, defined as the minimum  
 207 radius of a circle covering the whole graph minus one:  $N(r^{\max} + 1) = 1$ . The maximal radius and the *diameter*  
 208  $d^{\max}$  of a graph are also related, because  $r^{\max} \leq d^{\max} \leq 2r^{\max}$ . From these relations we can conclude the  
 209 following.



210 **Lemma 12.** For self-similar graphs (where  $N(r) \sim r^{-d}$ ), the diameter is  $d^{\max} \approx n^{1/d}$ , where  $d$  is the fractal  
 211 dimension.

212 In graphs where  $N(r) \sim e^{-\beta r}$ , the diameter is  $d^{\max} \approx \frac{\log n}{\beta}$ .

213 **Proof.** The diameter of a graph and the maximal radius are related as  $r^{\max} \leq d^{\max} \leq 2r^{\max}$ . Notice that, by  
 214 definition of the function  $N(r)$ , we have  $N(1) = n$ , where  $n$  is the number of nodes, and  $N(r^{\max} + 1) = 1$ .

215 Assuming  $N(r) = C r^{-d}$  and replacing  $r$  by 1 we get  $C = n$ . Then, replacing  $r$  by  $r^{\max} + 1$ , we get  
 216  $1 = N(r^{\max} + 1) = n (r^{\max} + 1)^{-d}$ . Hence,  $r^{\max} = n^{1/d} - 1$ .

217 Assuming  $N(r) = C e^{-\beta r}$  and replacing  $r$  by 1 we get  $C = n e^{\beta}$ . Then, replacing  $r$  by  $r^{\max} + 1$ , we get  
 218  $1 = N(r^{\max} + 1) = n e^{-\beta(r^{\max} + 1)}$ . Hence,  $r^{\max} = \frac{\log n}{\beta}$ .  $\square$

219 The diameter, as well as the typical distance<sup>5</sup>  $L$  of a graph, have been extensively used for characterizing  
 220 graphs. For instance, *small world graphs* [28] are characterized by small typical distance  $L \sim \log n$  and a large  
 221 clustering coefficient.

222 Notice that this characterization works well for *families* of graphs since these features can be expressed as a  
 223 function on the number of nodes. On the other hand, it may be inaccurate in the case of a particular graph, since  
 224 it is hard to assess whether its typical distance is “small” or if its clustering coefficient is “large”. Moreover, both  
 225 the diameter and the typical distance of a graph are hard-to-compute features. Notice that even though there is a  
 226 quadratic algorithm, computing them in huge graphs is a time-consuming task. This is the case of the application  
 227 SAT instances. In fact, our approximation to the fractal dimension can be computed more efficiently than the  
 228 diameter.

229 Additionally, the fractal dimension is independent of the graph size, and thus is a better feature to characterize  
 230 the instance structure. As a consequence, we can analyze formulas of the same family and very distinct size, all  
 231 having a similar structure. In particular, the fractal dimension will be similar in all the instances of the family,  
 232 and the shapes of the function  $N(r)$  will be resembling.

#### 233 4. The Self-Similar Structure of SAT Instances

234 In this paper, we have used the set of industrial formulas of the SAT Competition 2013<sup>6</sup>. They are 300  
 235 instances grouped into 19 families: *2d-strip-packing*, *bio*, *crypto-aes*, *crypto-des*, *crypto-gos*, *crypto-md5*,  
 236 *crypto-sha*, *crypto-vmc*, *diagnosis*, *hardware-bmc*, *hardware-bmc-ibm*, *hardware-cec*, *hardware-velev*, *planning*,  
 237 *scheduling*, *scheduling-pesp*, *software-bit-verif*, *software-bmc* and *termination*. In Table 1, we report some  
 238 statistics about these benchmarks. All these instances are *industrial*, in the sense that they come from a real-world  
 239 problem. During the paper, we compare them to random 3-CNF formulas.

240 We have conducted an exhaustive analysis of these 300 industrial SAT instances, and 90 random 3-CNF  
 241 formulas of  $10^5$  variables generated at different clause/variable ratios. We will see that most industrial instances  
 242 are self-similar and have a small fractal dimension, i.e.  $N(r) \sim r^{-d}$ , for small  $d$ . In random instances  $N(r)$   
 243 decays exponentially, i.e.  $N(r) \sim e^{-\beta r}$ . All experiments were carried out in a cluster of 9 nodes IBM dx360  
 244 M2, each of them with 32GB of RAM and 2 processors Intel(R) Xeon(R) CPU L5520 2.27 GHz, limiting all  
 245 experiments to a single core and to a maximum of 4GB of RAM.

246 Before presenting the results of this experimental analysis, we present a comparison of the BND and the  
 247 MEMB [27] algorithms, evaluating their performance and their accuracy. This analysis allows us to justify the  
 248 use of the BND algorithm to compute the fractal dimension of SAT instances.

249 We execute the BND and the MEMB algorithms for the set of 300 industrial instances of the SAT  
 250 Competition 2013. In this experiment, we set a timeout limit of half an hour. The BND algorithm is able  
 251 to complete the computation in all the 300 instances. On the contrary, the MEMB algorithm only finishes in 17  
 252 instances. Also, the BND algorithm spends an average run-time of 0.11 seconds, while the MEMB algorithms  
 253 spends 607.2 seconds in average. Finally, the function  $N^C(r)$  approximated by both algorithms is very alike.

<sup>5</sup> The typical distance of a graph is the average of the distances between any two nodes.

<sup>6</sup> <http://satcompetition.org/2013/>

Family	#instances	avg. #vars.	avg. #clauses
2d-strip-packing	5	12753.20	1267404.40
bio	5	46376.40	471256.00
crypto-aes	11	25285.90	82980.81
crypto-des	9	31590.22	95015.66
crypto-gos	30	1892.93	22883.60
crypto-md5	11	66894.00	267579.54
crypto-sha	30	4773.33	151027.86
crypto-vmcpc	8	1061.50	171947.00
diagnosis	26	235279.42	1129797.46
hardware-bmc-ibm	4	1335863.00	5234846.00
hardware-bmc	3	114448.33	341918.00
hardware-cec	30	255757.16	754358.03
hardware-velev	21	282973.85	7478281.09
planning	25	469285.60	3323685.24
scheduling-pesp	30	37454.30	270133.10
scheduling	30	160310.13	722325.73
software-bit-verif	14	131492.71	378409.78
software-bmc	3	11216817.00	32697150.00
termination	5	219846.80	950546.60
SATComp2013	300	234984.90	1469306.71

**Table 1.** Statistics about the application benchmarks of the SAT Competition 2013.

254 This can be seen in Fig. 2, where it is represented the obtained approximation in the 17 SAT formulas solved  
 255 by both algorithms in the given timeout. In the rest of our experimental analysis, we will only use the BND  
 256 algorithm, given that the shapes of both approximations are very similar and that the MEMB algorithm is much  
 257 slower.

258 We also emphasize that the MEMB algorithm is slightly more accurate in some cases. This explains why a  
 259 few upper bounds of  $N^C(r)$  that BND calculates are above the ones calculated by MEMB. The exact values of  
 260  $N^C(r)$  are probably even lower, especially for high values of  $r$  (where the approximation is less accurate).

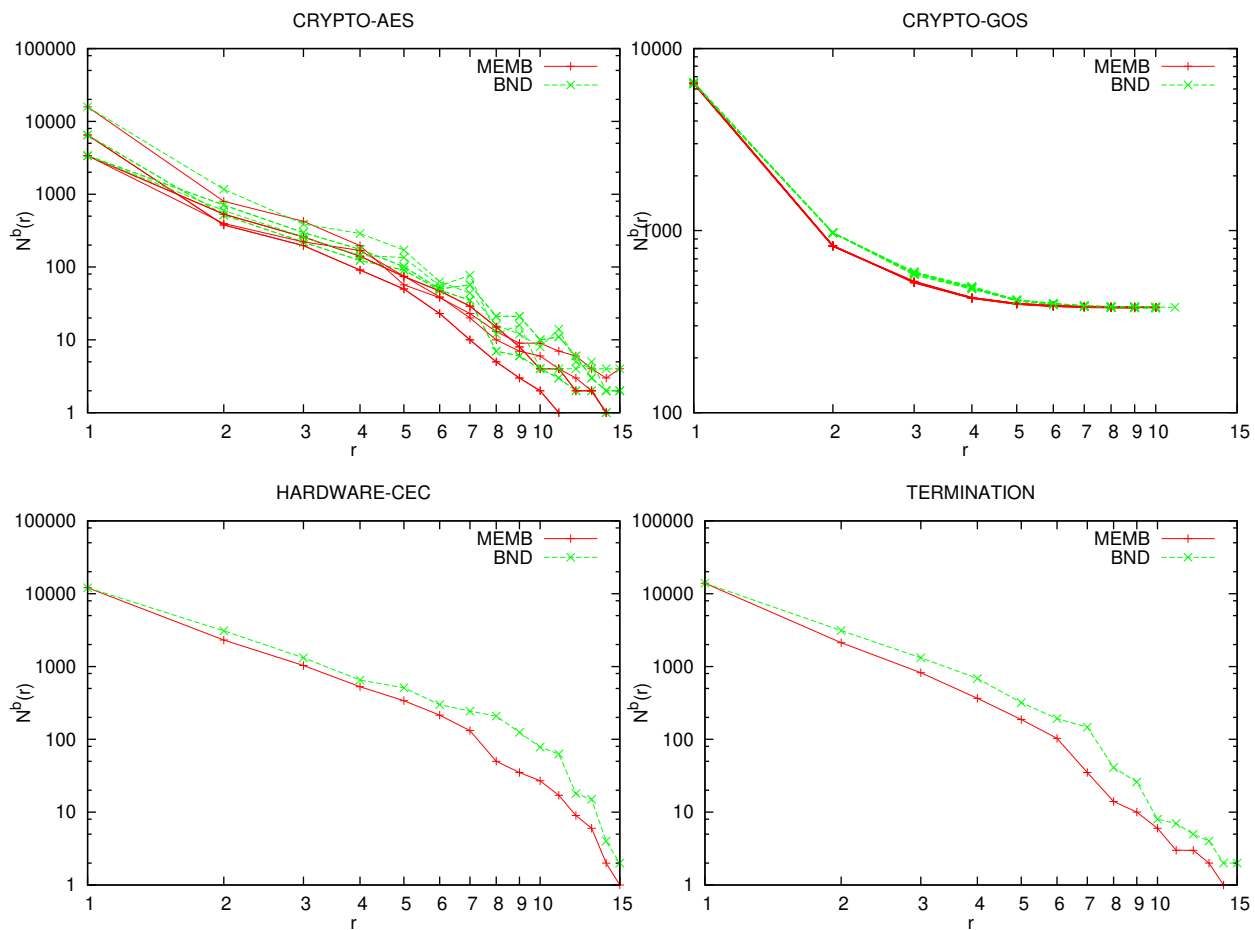
261 The VIG model that represent a random 2-SAT instance is exactly an Erdős-Rényi graphs. These formulas  
 262 have a satisfiability threshold at  $m/n = 1$ , commonly known as the phase transition point. This means that all  
 263 formula below (resp. above) this ratio are satisfiable (resp. unsatisfiable). It is also known that at  $m/n = 0.5$   
 264 there is a percolation threshold. This means that the VIG of formulas above this ratio have a major connected  
 265 component, whereas below this threshold the VIG graph breaks into many connected components. In the  
 266 percolation point the formula is self-similar with a fractal dimension  $d = 2$ . Above this ratio, the function  $N^V(r)$   
 267 decays exponentially. To the best of our knowledge, there is no known result of this kind for random 3-CNF  
 268 formulas.

269 In Fig. 3 we represent the functions  $N^V(r)$  and  $N^C(r)$ , approximated by the BND algorithm, for families  
 270 of random 3-CNF instances with different clause/variable ratios. Although it is not showed in this figure, we  
 271 observed that these functions does not depend on the number of variables of the family; they only depend on  
 272 their ratio of clause-to-variable  $m/n$ .

273 In the phase transition point  $m/n = 4.25$ , the function  $N^V(r)$  exhibit an exponential decay with  $\beta = 2.3$ ,  
 274 i.e.,  $N^V(r) \sim e^{-2.3r}$ . Hence,  $r^{max} = \frac{\log n}{2.3} + 1$ . For example, with  $n = 10^5$  variables random SAT formulas  
 275 have a radius  $r^{max} \approx 6$ . The decay factor is even higher for higher values of  $m/n$ . The same phenomenon is  
 276 observed in the CVIG model. However, the exponential decay is  $\beta = 1.16 \approx 2.3/2$  in the family of formulas at  
 277 the phase transition point. Hence, the decay in the VIG model is approximately double of the decay in the CVIG,  
 278 as expected by Lemma 11. A percolation threshold is experimentally found at  $m/n \approx 0.17$ , where the main  
 279 connected component also exhibits a fractal dimension  $d = 2$ .

280 In Fig. 4 we represent the function  $N^C(r)$  for some application or industrial families of SAT instances. It  
 281 can be observed that most of these families exhibit self-similar structure, with most fractal dimensions ranging  
 282 between 2 and 4. All the formulas of the same family exhibit the same fractal dimension in many of the families,  
 283 where this dimension seems to be a characteristic feature of the family. In many of the industrial families, all





**Figure 2.** Upper bounds for  $N^b(r)$  obtained with MEMB and BND algorithms, for the 17 industrial instances that MEMB is able to compute in 30 minutes, grouped by families.

284 instances have the same fractal dimension, being this dimension a particular property of the family. See, for  
 285 instance, families *crypto-sha* or *diagnosis*. It is interesting to see that the formula size does not alter the value of  
 286 the fractal dimension (in our representation in logarithmic scale, they all share the same function shape).

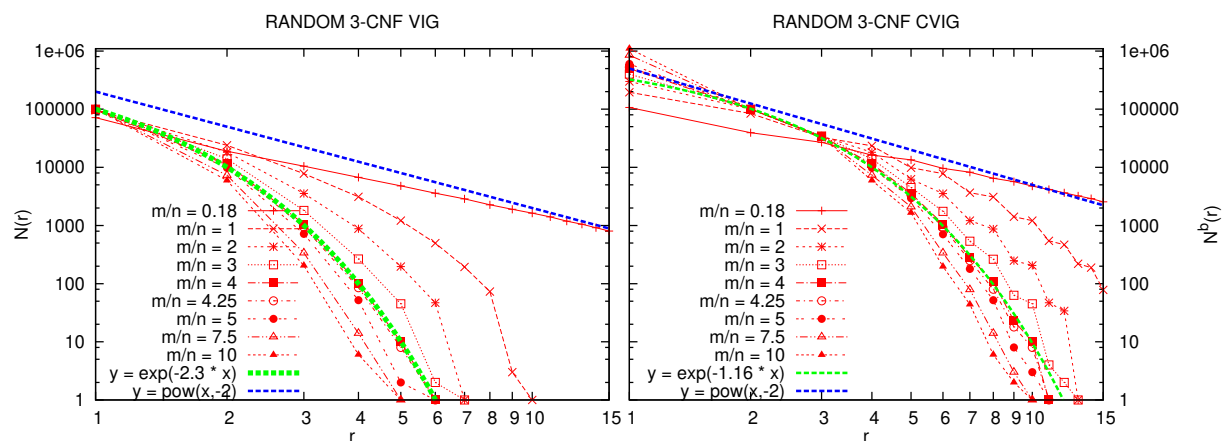
287 In general, we observe that for small values of  $r$ , the polynomial decay is clearer, with an almost identical  
 288 slope for all the formulas of the family.

289 On the other hand, for big values of  $r$ , we must make some remarks.

290 First, the upper bound on  $N^C(r)$  that we compute may be inaccurate. Second, we identify two events. First,  
 291 in some cases although the whole function can not be approximated by an exponential function, there is an abrupt  
 292 decay (see some *hardware-cec* or *termination* instances, for example). This can be due to a small set of edges  
 293 that connect distant part of the graph. Notice they may drastically reduce the number of circles required to cover  
 294 the graph when the value of  $r$  is high, but they would have no effect for smaller values. Second, there is a long  
 295 tail in some other cases (see *hardware-bmc-ibm*, for instance). This is due to the existence of (small) unconnected  
 296 components in the graph. However, this tail disappears if we compute  $N(r)$  only for the major component.<sup>7</sup>

297 Finally, there are some families where the  $N(r)$  function decays exponentially, i.e., they are not self-similar.  
 298 This is the case, for instance, of the *hardware-velev* family. However, we emphasize that this kind of families are  
 299 rare, since most of the industrial formulas and families of our analysis do have self-similar structure, i.e., with a  
 300  $N(r)$  function decaying polynomially.

<sup>7</sup> We can subtract from  $N(r)$  the number of unconnected components, as an approximation, since most are covered with a few circles.



**Figure 3.** Functions  $N^V(r)$  for VIG (left), and  $N^C(r)$  for CVIG (right), for 3-CNF random formulas with distinct values of  $m/n$ . Formulas are generated using  $n = 10^5$  variables and taking the major connected component, except for  $m/n = 0.18$ , where  $n = 10^6$ .

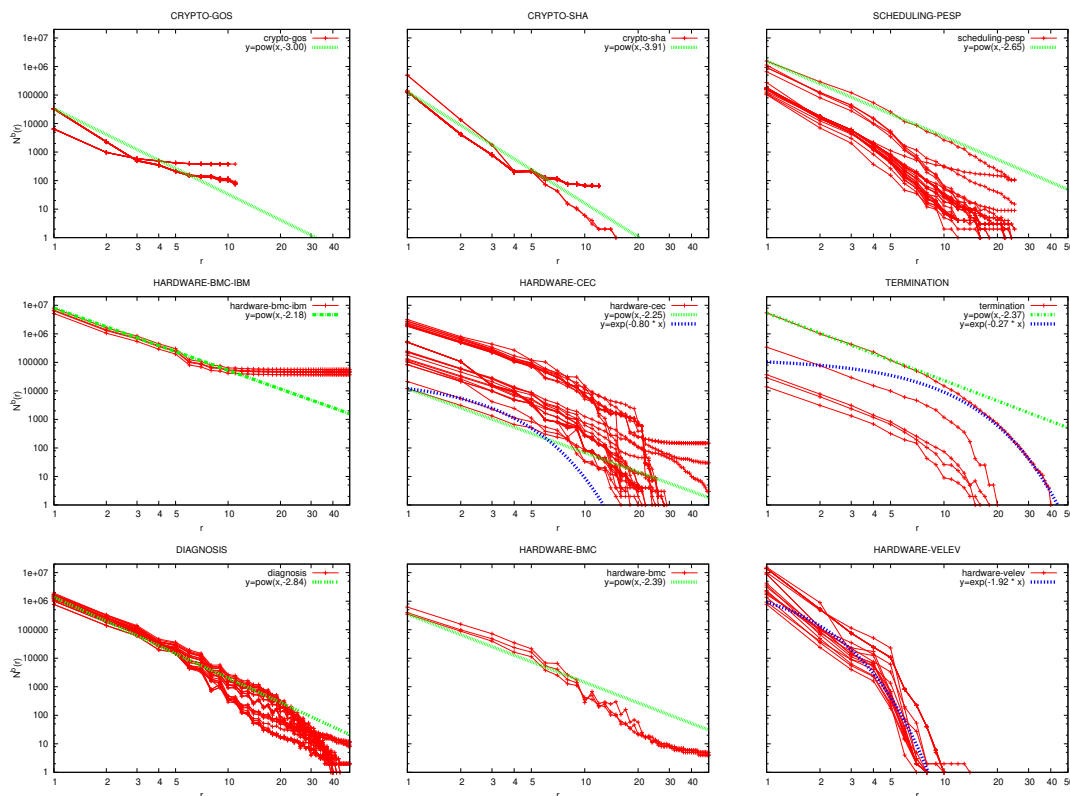
#### 4.1. Fractal Dimension at Fine-Grained Scale

If a graph exhibits self-similar structure, the same structure can be observed at all scales. This means that we could group a subset of nodes covered (i.e., the ones covered by a circle) into a single node, without altering the same structure in the resulting graph. In our evaluation, we observe that this happens when  $r$  is small (the function behavior is  $N(r) \approx C r^{-d}$  when  $r$  is small). However, this does not always happen for big values of  $r$ . This can be because the self-similar structure disappears at a coarse-grained scale, or because the function  $N(r)$  is not well approximated. When a SAT instance has a small unsatisfiability proof, a small cycle can be observed in the graph models representing the formula (e.g., VIG and CVIG), suggesting that the fine-grained scale of the fractal dimension is more relevant. In other words, we conjecture that is more important the value of the fractal dimension at fine-grained scale, i.e. the slope of the function  $N(r)$  for small values of  $r$ , rather than the existence of a self-similar structure. Therefore, in our evaluation, we note these fine-grained dimensions as  $d^V$  and  $d^C$  for the VIG and CVIG, respectively. We compute them as the interpolation, by linear regression, of  $\log N(r)$  vs.  $\log r$ . We use the values of  $N^V(r)$  and  $N^C(r)$ , for  $r = 1, \dots, 6$ . Experimentally, we see that these approximations are accurate enough.

### 5. The Self-Similar Structure during SAT Solver Search

CDCL SAT solvers add during their execution new clauses learned from the conflicts the solver finds. A conflict represents a (partial) assignment of the variables which leads to a contradiction (i.e., an unsatisfied clause). These conflicts can be learned in order to prevent the solver to find them in the future. When a learned clause is unary (i.e., it only has one literal), this literal can be propagated (i.e., assigned), and this results in a simplification of the given formula. Given a unary clause  $x$ , clauses with literal  $x$  are completely removed (i.e., they are satisfied by  $x$ ), and literals  $\neg x$  are removed from the formula. On the other hand, learned clauses of bigger length may create new relations between variables, i.e., new edges in the graph representations of the formula.

Both, the simplification caused by learning unary clauses, and the addition of longer learned clauses, can alter the fractal dimension of the SAT instance. In general graphs, the addition of edges (preserving the nodes) can only increase its dimension, because existing circles may cover more nodes (due to the new edges), and thus the number  $N(r)$  of circles required to cover the graph may decrease, whereas  $N(1)$  (the number of nodes) remains unchanged. As a consequence, this can only increase the fractal dimension since the slope of  $N(r)$  increases. In the VIG model, learning non-unary clauses only introduces new edges, thus as we argued the fractal dimension can only increase. On the contrary, this argument is more complicated in the CVIG model. Let  $N(r)$  be the original number of circles and  $N'(r)$  the minimum number of circles needed after adding  $L$  learned clauses. We have  $N'(1) = N(1) + L$ , since we add  $L$  new nodes. For  $r > 1$ , the whole graph can be covered with the



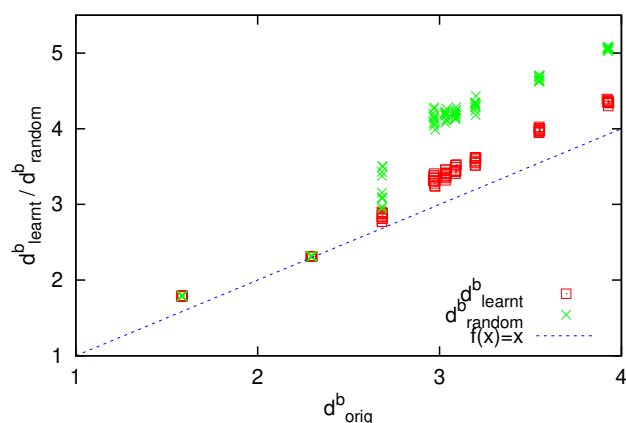
**Figure 4.** Function  $N^C(r)$  for some industrial SAT formulas grouped by families.

333 former circles and using  $L$  new circles, used to cover the new nodes, thus  $N'(r) \leq N(r) + L$  in the worst case.  
 334 Therefore, the dimension could only increase. However, CDCL SAT solvers do not add unary clauses to the  
 335 formula, but propagate them. And this decreases the value of  $N(1)$  in both models, while increasing the value of  
 336  $N(r)$  (with  $r > 1$ ) because of the elimination and/or simplification of the clauses after the propagation(s). As a  
 337 result, the fractal dimension can only decrease. In the previous argument, we assume that  $\log N(r)$  is perfectly  
 338 linear on  $\log r$ , and  $N(r)$  is computed exactly. In practice, we compute an approximation of  $N(r)$ . Moreover, we  
 339 compute the dimension by linear regression, since points are not aligned. Hence, the variation in the dimension  
 340 due to learning is rather unpredictable.

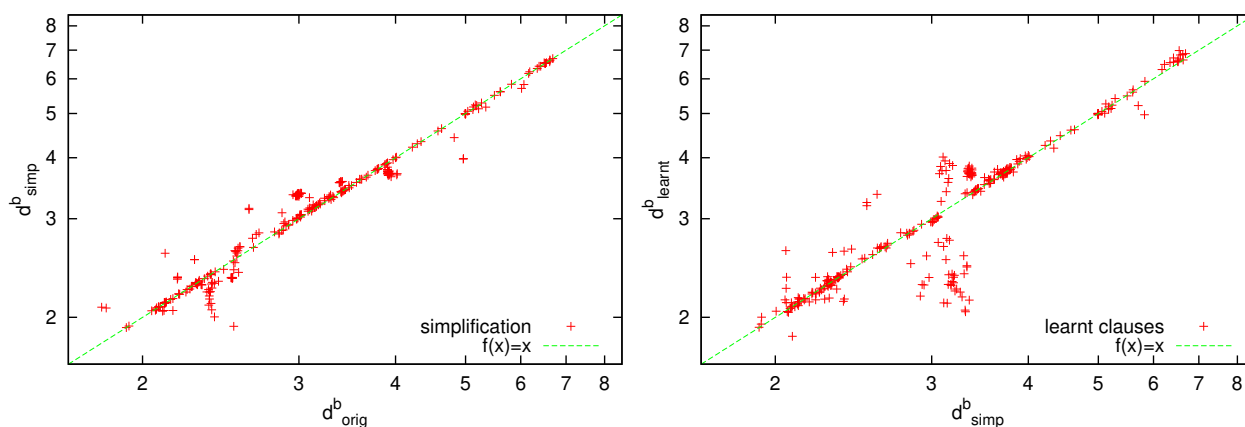
341 In order to study the evolution of the fractal dimension during SAT solver search, we have conducted  
 342 the following experiments. First, we analyze this phenomenon in random 3-CNF instances with  $10^5$  variables  
 343 and distinct clause/variable ratios. It is very unlikely that solving these formulas produces any unary learned  
 344 clauses, hence neither nodes nor edges are removed. In Fig. 5, we compare the original dimension  $d_{orig}^C$  of the  
 345 formulas and its the fractal dimension  $d_{learned}^C$  after adding learned clauses. It can be observed that learning  
 346 new clauses increases the dimension of the formula, as expected theoretically, being is bigger for formulas with  
 347 higher clause/variable ratio. This increase can be *quantified* repeating the same experiment and replacing the  
 348 clauses learned by the solver by randomly generated clauses of the same length, and computing the new fractal  
 349 dimension  $d_{random}^C$ ; the results are also shown in Fig. 5. In this second experiment it can be observed that the  
 350 increase in the dimension is bigger when learning random clauses:  $d_{random}^C \geq d_{learned}^C \geq d_{orig}^C$ . This means that  
 351 actual learned clauses, even in random  $k$ -CNF formulas, tend to connect variables that were already *close* in the  
 352 graph representation of the formula.

353 In application SAT formulas, in contrast, the solver usually learns some unary clauses. In Fig. 6, we  
 354 represent separately the effects of simplifying the formula after *learning* unary clauses (left), and the effect of  
 355 learning longer clauses (right), after  $10^3$  conflicts.

356 In most of the application benchmarks, simplifying the formulas with unary clauses slightly in general  
 357 decreases its fractal dimension, as expected. This dimension is represented as  $d_{simp}^C$ . In fact, in some instances



**Figure 5.** Relation between the original fractal dimension  $d_{orig}^C$ , and the dimension  $d_{learnt}^C$  after adding learned clauses, or after adding random clauses  $d_{rand}^C$ , in random 3-CNF formulas. learned clauses are computed after  $10^3$  conflicts.

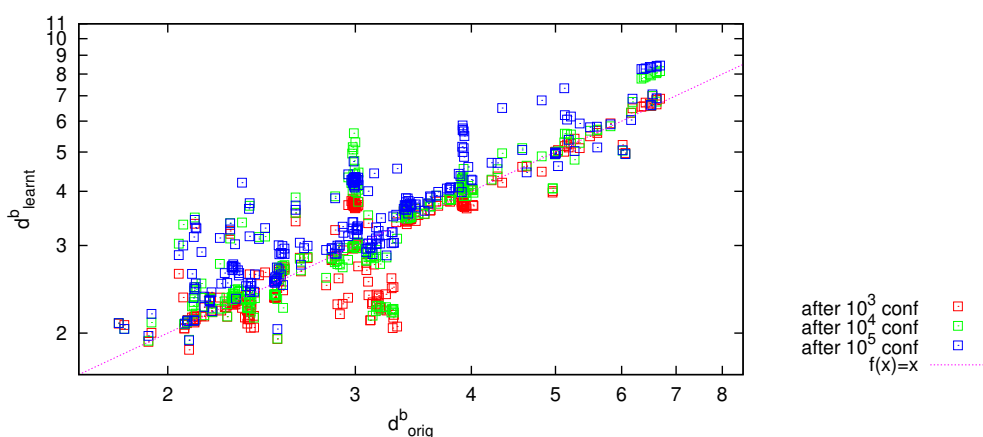


**Figure 6.** Relation between the original fractal dimension  $d_{orig}^C$  and the fractal dimension  $d_{simp}^C$  after simplifying the formula with the unary learned clauses (left), and relation between the fractal dimension  $d_{simp}^C$  and the fractal dimension  $d_{learnt}^C$  after simplification and adding learned clauses (right), for all industrial formulas. learned clauses are the result of  $10^3$  conflicts.

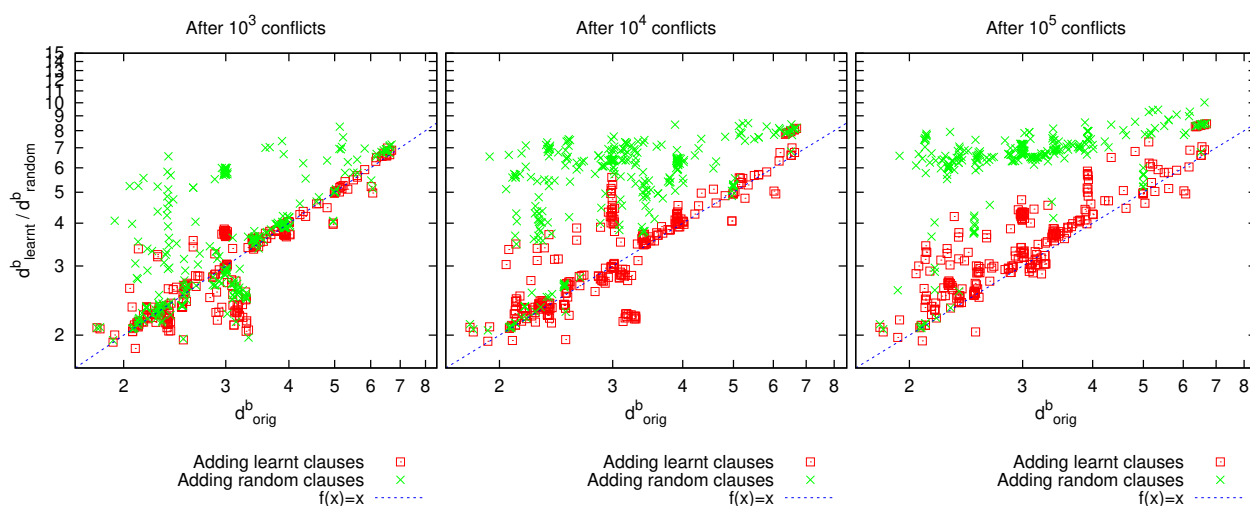
358 there is a huge difference between  $d_{simp}^C$  and  $d_{orig}^C$ . This is the case of families *crypto-md5* and *crypto-sha*.  
 359 In these benchmarks, propagating unary clauses is more important, and produces more changes in the fractal  
 360 dimension. Nevertheless, in a few instances the fractal dimension  $d_{simp}^C$  slightly increases. Analyzing the effects  
 361 of learning longer clauses, it can be observed in most of the instances that the dimension  $d_{learnt}^C$  after adding  
 362 such learned clauses slightly increases the fractal dimension w.r.t.  $d_{simp}^C$ , as expected.

363 This increase is particularly remarkable in the families *crypto-aes*, *crypto-gos* and *crypto-vmc*. This  
 364 suggests that learned clauses connect distant part of the formula, hence  $N^C(r)$  drops off. Again, we also observe  
 365 some cases where the dimension  $d_{learnt}^C$  decreases due to clause learning.

366 In Fig. 7, we represent the variations in the fractal dimension after  $10^3$ ,  $10^4$  and  $10^5$  conflicts. It can be  
 367 observed that fractal dimensions can increase or decrease slightly after  $10^3$  conflicts, depending on the impact  
 368 of simplifications caused by unary clauses (producing a lower fractal dimension) or adding learned clauses  
 369 (producing a higher the fractal dimension). Nevertheless, after a longer number of conflicts (e.g.,  $10^5$  conflicts),  
 370 there is a more clear tendency: the fractal dimension increases in most of the cases. This results matches the  
 371 expected behavior of the solver, where most of the unary clauses are learned in the early stages of the search, and  
 372 explains that the dimensions tends to increase.



**Figure 7.** Relation between the original fractal dimension  $d_{orig}^C$  and the fractal dimensions after learning clauses  $d_{learned}^C$ , in industrial formulas.



**Figure 8.** Relation between the original fractal dimension  $d_{orig}^C$  and the fractal dimensions after adding learned clauses  $d_{learned}^C$ , or after adding random clauses  $d_{random}^C$ , in industrial formulas.

373 Finally, we analyze whether learned clauses (tend to) connect distant part of the formula in these applications  
 374 benchmarks. For this purpose, we compare the fractal dimension  $d_{learned}^C$  of the formulas after adding learned  
 375 clauses and the dimension  $d_{random}^C$  where the previous learned clauses are replaced by randomly generated  
 376 clauses of the same length. In Fig. 8, we represent these two dimensions. It can be observed that  $d_{random}^C$  is  
 377 considerably higher than  $d_{learned}^C$ .

378 This means that the clauses learned by the solver do not contribute to decrease  $N(r)$ , i.e., the number of  
 379 circles required to cover the graph is similar. On the contrary, these random clauses would significantly decrease  
 380  $N(r)$ . Therefore, learned clauses mainly connect variables (nodes) that were already covered by the same circles,  
 381 i.e., nodes already connected or close in the graph, and thus, clause learning acts quite locally in the formula.

## 382 6. Related Work

383 In the last ten years, there have been many works on the analysis of real-world SAT instances, seen as  
 384 graphs [22–24,29–33]. The main aim of these studies is to better understand the success of CDCL techniques  
 385 and help to improve them. They can also help develop new benchmarks generation models with more realistic  
 386 properties [34–40]. This problem has been identified as an important challenge in the SAT community [41–43].

387 The definition of fractal dimension for graphs is based on the notion of *box covering* [26].  
 388 In [27] it is described a box-covering algorithm particularly suitable for complex networks, called



389 Maximum-Excluded-Mass-Burning (MEMB). They have used it to demonstrate the existence of self-similarity  
 390 in many real-world networks [44,45]. There are other definitions of dimension for complex networks, like the  
 391 one based on the *mass* (average number of vertices in a box) of the clusters [46]. The algorithm we propose  
 392 in this paper is a variation of the MEMB algorithm, with the addition of a heuristic that allows us to process  
 393 huge graphs, like the ones obtained from industrial SAT benchmarks. There are several models for generating  
 394 scale-free graphs with fractal dimension [47–49].

## 395 7. Conclusions

396 The community on *complex networks* have developed methods and models to better understand the structure  
 397 of real-world graphs. The classical Erdős-Rényi model, the first model in graph theory, cannot be used for  
 398 the study of *real-world* networks, since this model does not generate the *structural* features these networks  
 399 exhibit. Similarly, classical random SAT formulas cannot be used to better understand the underlying structure  
 400 of industrial SAT instances and its relation to SAT solver performance. These instances are characterized by a  
 401 particular structure, which may explain their distinct nature with respect to random formulas.

402 In this work, we have analyzed the *self-similarity* (or *fractal dimension*) of industrial SAT formulas, and  
 403 how this structure evolves during the SAT solver search. We have found that most of industrial formulas show a  
 404 self-similar structure, while random formulas do not. Their fractal dimension ranges, in general, between 2 and 4.  
 405 Fractal dimension, typical distances and graph diameter are related (small dimension implies big distance and  
 406 big diameter). Hence, industrial SAT instances have a big diameter. Intuitively, that means that we need long  
 407 chains of implications to propagate a variable instantiation to others. Moreover, we have observed that fractal  
 408 dimension increases due to learned clauses. However, this increase is bigger if we substitute learned clauses by  
 409 random clauses of the same size. Therefore, learning *does not* contribute to connect distant parts of the formula  
 410 very much. In other words, it works *locally*.

411 We think that the present study provides a step towards a practical explanation of why some SAT solvers  
 412 perform better on industrial instances, and they do not on random formulas.

413 **Funding:** This work is partially supported by the EU H2020 Research and Innovation Programme under the LOGISTAR  
 414 project (Grant Agreement No. 769142), MINECO-FEDER projects RASO (TIN2015-71799-C2-1-P) and TASSAT3  
 415 (TIN2016-76573-C2-2-P), the Spanish Ministerio de Economía y Competitividad under the EXASOCO project (ref.  
 416 PGC2018-101216-B-I00), including European Regional Development Funds (ERDF). The third author is also supported by a  
 417 MICINN Juan de la Cierva fellowship (ref. FJCI-2017-32420).

## 418 Abbreviations

419 The following abbreviations are used in this manuscript:

420 SAT	the Boolean satisfiability problem
CDCL	Conflict-Driven Clause Learning
DPLL	Davis-Putnam-Logemann-Loveland algorithm
VSIDS	Variable State Independent Decaying Sum heuristics
421 VIG	Variable Incidence Graph
CVIG	Clause-Variable Incidence Graph
MEMB	Maximum-Excluded-Mass-Burning algorithm
BND	Burning by Node Degree algorithm

422

- 423 1. Biere, A. Bounded Model Checking. In *Handbook of Satisfiability*; IOS Press, 2009; pp. 457–481.
- 424 2. Rintanen, J. Planning and SAT. In *Handbook of Satisfiability*; IOS Press, 2009; pp. 483–504.
- 425 3. Kroening, D. Software Verification. In *Handbook of Satisfiability*; IOS Press, 2009; pp. 505–532.
- 426 4. Silva, J.P.M.; Lynce, I.; Malik, S. Conflict-Driven Clause Learning SAT Solvers. In *Handbook of Satisfiability*; IOS  
 427 Press, 2009; pp. 131–153.
- 428 5. Davis, M.; Putnam, H. A Computing Procedure for Quantification Theory. *J. ACM* **1960**, *7*, 201–215.  
 429 doi:10.1145/321033.321034.

- 430 6. Davis, M.; Logemann, G.; Loveland, D.W. A machine program for theorem-proving. *Commun. ACM* **1962**,  
431 5, 394–397. doi:10.1145/368273.368557.
- 432 7. Silva, J.P.M.; Sakallah, K.A. GRASP: A Search Algorithm for Propositional Satisfiability. *IEEE Trans. Computers*  
433 **1999**, 48, 506–521. doi:10.1109/12.769433.
- 434 8. Moskewicz, M.W.; Madigan, C.F.; Zhao, Y.; Zhang, L.; Malik, S. Chaff: Engineering an Efficient SAT Solver. Proc.  
435 of the 38th Design Automation Conf. (DAC'01), 2001, pp. 530–535.
- 436 9. Audemard, G.; Simon, L. Predicting Learnt Clauses Quality in Modern SAT Solvers. Proc. of the 21st Int. Joint  
437 Conf. on Artificial Intelligence (IJCAI'09), 2009, pp. 399–404.
- 438 10. Gomes, C.P.; Selman, B.; Kautz, H.A. Boosting Combinatorial Search Through Randomization. Proc. of the 15th  
439 Nat. Conf. on Artificial Intelligence (AAAI'98), 1998, pp. 431–437.
- 440 11. Eén, N.; Biere, A. Effective Preprocessing in SAT Through Variable and Clause Elimination. Proc. of the 8th Int.  
441 Conf. on Theory and Applications of Satisfiability Testing (SAT'05), 2005, pp. 61–75.
- 442 12. Sakallah, K.A.; Marques-Silva, J. Anatomy and Empirical Evaluation of Modern SAT Solvers. *Bulletin of the EATCS*  
443 **2011**, 103, 96–121.
- 444 13. Elffers, J.; Giráldez-Cru, J.; Gocht, S.; Nordström, J.; Simon, L. Seeking Practical CDCL Insights from Theoretical  
445 SAT Benchmarks. Proc. of the 27th Int. Joint Conf. on Artificial Intelligence (IJCAI'18), 2018, pp. 1300–1308.  
446 doi:10.24963/ijcai.2018/181.
- 447 14. Liang, J.H.; Ganesh, V.; Poupart, P.; Czarnecki, K. Learning Rate Based Branching Heuristic for SAT Solvers. Proc.  
448 of the 19th Int. Conf. on Theory and Applications of Satisfiability Testing (SAT'16), 2016, pp. 123–140.
- 449 15. Biere, A. CaDiCaL at the SAT Race 2019. Proceedings of SAT Race 2019 : Solver and Benchmark Descriptions,  
450 2019, pp. 8–9.
- 451 16. Williams, R.; Gomes, C.P.; Selman, B. Backdoors To Typical Case Complexity. Proc. of the 18th Int. Joint Conf. on  
452 Artificial Intelligence (IJCAI'01), 2003, pp. 1173–1178.
- 453 17. Gomes, C.P.; Selman, B. Problem Structure in the Presence of Perturbations. Proc. of the 14th Nat. Conf. on  
454 Artificial Intelligence (AAAI'97), 1997, pp. 221–226.
- 455 18. Hogg, T. Refining the Phase Transition in Combinatorial Search. *Artif. Intell.* **1996**, 81, 127–154.
- 456 19. Gent, I.P.; Hoos, H.H.; Prosser, P.; Walsh, T. Morphing: Combining Structure and Randomness. Proc. of the 16th  
457 Nat. Conf. on Artificial Intelligence (AAAI'99), 1999, pp. 654–660.
- 458 20. Järvisalo, M.; Niemelä, I. The effect of structural branching on the efficiency of clause learning SAT solving: An  
459 experimental study. *J. Algorithms* **2008**, 63, 90–113.
- 460 21. Erdős, P.; Rényi, A. On Random Graphs. *Publicationes Mathematicae* **1959**, 6, 290–297.
- 461 22. Ansótegui, C.; Giráldez-Cru, J.; Levy, J. The Community Structure of SAT Formulas. Proc. of the 15th Int. Conf. on  
462 Theory and Applications of Satisfiability Testing (SAT'12), 2012, pp. 410–423.
- 463 23. Ansótegui, C.; Bonet, M.L.; Levy, J. On the Structure of Industrial SAT Instances. Proc. of the 15th Int. Conf. on  
464 Principles and Practice of Constraint Programming (CP'09), 2009, pp. 127–141.
- 465 24. Ansótegui, C.; Bonet, M.L.; Giráldez-Cru, J.; Levy, J. The Fractal Dimension of SAT Formulas. Proc. of the 7th Int.  
466 Joint Conf. on Automated Reasoning (IJCAR'14), 2014, pp. 107–121.
- 467 25. Giráldez-Cru, J. Beyond the Structure of SAT Formulas. PhD thesis, Universitat Autònoma de Barcelona, 2016.
- 468 26. Mandelbrot, B.B. *The fractal geometry of nature*; Macmillan, 1983.
- 469 27. Song, C.; Gallos, L.K.; Havlin, S.; Makse, H.A. How to calculate the fractal dimension of a complex network: the  
470 box covering algorithm. *Journal of Statistical Mechanics: Theory and Experiment* **2007**, 2007, P03006.
- 471 28. Walsh, T. Search in a Small World. Proc. of the 16th Int. Joint Conf. on Artificial Intelligence (IJCAI'99), 1999, pp.  
472 1172–1177.
- 473 29. Katsirelos, G.; Simon, L. Eigenvector Centrality in Industrial SAT Instances. Proc. of the 19th Int. Conf. on  
474 Principles and Practice of Constraint Programming (CP'12), 2012, pp. 348–356.
- 475 30. Ansótegui, C.; Giráldez-Cru, J.; Levy, J.; Simon, L. Using community structure to detect relevant learnt clauses. Proc.  
476 of the 18th Int. Conf. on Theory and Applications of Satisfiability Testing (SAT'15), 2015, pp. 238–254.
- 477 31. Ansótegui, C.; Bonet, M.L.; Giráldez-Cru, J.; Levy, J. On the Classification of Industrial SAT Families.  
478 Proc. of the 18th Int. Conf. of the Catalan Association for AI (CCIA'15), 2015, pp. 163–172.  
479 doi:10.3233/978-1-61499-578-4-163.
- 480 32. Baud-Berthier, G.; Giráldez-Cru, J.; Simon, L. On the Community Structure of Bounded Model Checking SAT  
481 Problems. Proc. of the 20th Int. Conf. on Theory and Applications of Satisfiability Testing (SAT'17), 2017, pp.  
482 65–82. doi:10.1007/978-3-319-66263-3\_5.

- 483 33. Ansótegui, C.; Bonet, M.L.; Giráldez-Cru, J.; Levy, J. Structure features for SAT instances classification. *J. Applied*  
484 *Logic* **2017**, *23*, 27–39. doi:10.1016/j.jal.2016.11.004.
- 485 34. Slater, A. Modelling More Realistic SAT Problems. Proc. of the 15th Australian Joint Conf. on Artificial Intelligence  
486 (AJCAI'02), 2002, pp. 591–602.
- 487 35. Ansótegui, C.; Bonet, M.L.; Levy, J. Towards Industrial-Like Random SAT Instances. Proc. of the 21st Int. Joint  
488 Conf. on Artificial Intelligence (IJCAI'09), 2009, pp. 387–392.
- 489 36. Burg, S.; Kaufmann, M.; Kottler, S. Creating Industrial-Like SAT Instances by clustering and reconstruction. Proc.  
490 of the 15th Int. Conf. on Theory and Applications of Satisfiability Testing (SAT'12), 2012, pp. 471–472.
- 491 37. Järvisalo, M.; Kaski, P.; Koivisto, M.; Korhonen, J.H. Finding Efficient Circuits for Ensemble Computation. Proc. of  
492 the 15th Int. Conf. on Theory and Applications of Satisfiability Testing (SAT'12), 2012, pp. 369–382.
- 493 38. Giráldez-Cru, J.; Levy, J. A modularity-based random SAT instances generator. Proc. of the 24th Int. Joint Conf. on  
494 Artificial Intelligence (IJCAI'15), 2015, pp. 1952–1958.
- 495 39. Giráldez-Cru, J.; Levy, J. Generating SAT instances with community structure. *Artif. Intell.* **2016**, *238*, 119–134.
- 496 40. Giráldez-Cru, J.; Levy, J. Locality in Random SAT Instances. Proc. of the 26th Int. Joint Conf. on Artificial  
497 Intelligence (IJCAI'17), 2017, pp. 638–644.
- 498 41. Selman, B.; Kautz, H.A.; McAllester, D.A. Ten Challenges in Propositional Reasoning and Search. Proc. of the 15th  
499 Int. Joint Conf. on Artificial Intelligence (IJCAI'97), 1997, pp. 50–54.
- 500 42. Kautz, H.A.; Selman, B. Ten Challenges Redux: Recent Progress in Propositional Reasoning and Search. Proc. of  
501 the 9th Int. Conf. on Principles and Practice of Constraint Programming (CP'03), 2003, pp. 1–18.
- 502 43. Dechter, R. *Constraint Processing*; Morgan Kaufmann, 2003.
- 503 44. Song, C.; Havlin, S.; Makse, H.A. Self-similarity of complex networks. *Nature* **2005**, *433*, 392–395.
- 504 45. Song, C.; Havlin, S.; Makse, H.A. Origins of fractality in the growth of complex networks. *Nature Physics* **2006**,  
505 *275*, 275.
- 506 46. Shanker, O. Defining Dimension of a Complex Network. *Modern Physics Letters B* **2007**, *21*, 321–326.
- 507 47. Dorogovtsev, S.N.; Goltsev, A.V.; Mendes, J.F.F. Pseudofractal scale-free web. *Physical Review E* **2002**, *65*, 066122.
- 508 48. Ravasz, E.; Barabási, A.L. Hierarchical organization in complex networks. *Phys. Rev. E* **2003**, *67*, 026112.  
509 doi:10.1103/PhysRevE.67.026112.
- 510 49. Jung, S.; Kim, S.; Kahng, B. Geometric fractal growth model for scale-free networks. *Phys. Rev. E* **2002**, *65*, 056101.  
511 doi:10.1103/PhysRevE.65.056101.