

# The essentials of the SAT 2003 competition

Daniel Le Berre<sup>1</sup> and Laurent Simon<sup>2</sup>

<sup>1</sup> CRIL, Université d'Artois,  
Rue Jean Souvraz SP 18 – F 62307 Lens Cedex, France  
[leberre@cril.univ-artois.fr](mailto:leberre@cril.univ-artois.fr)

<sup>2</sup> LRI, Université Paris-Sud  
Bâtiment 490, U.M.R. CNRS 8623 – 91405 Orsay Cedex, France  
[simon@lri.fr](mailto:simon@lri.fr)

**Abstract.** The SAT 2003 Competition ran in February – May 2003, in conjunction with SAT'03 (the Sixth Fifth International Symposium on the Theory and Applications of Satisfiability Testing). One year after the SAT 2002 competition, it was not clear that significant progress could be made in the area in such a little time. The competition was a success – 34 solvers and 993 benchmarks, needing 522 CPU days – with a number of brand new solvers. Several 2003 competitors were even able to solve within 15mn benchmarks remained unsolved within 6 hours by 2002 competitors. We report here the essential results of the competition, interpret and statistically analyse them, and at last provide some suggestions for the future competitions.

## 1 Introduction

Organizing a competition is a large effort, made by the whole community, both for submitters and organizers. It is thus essential to ensure that this effort is keeping worthwhile for the community. Clearly, such an event is keeping the community excited and allows a good visibility of the SAT field outside its own frontiers. For the second edition of the competition, the challenges were to bring up new solvers, to measure progress made since last year (positive if possible!) and, more important, to identify new (good) techniques and new challenges for SAT solvers.

The competition was a success in terms of submissions, with 17 benchmarks submitters (benchmarks or generators) and 21 solvers submitters. 34 solvers entered the contest and we selected 993 benchmarks, partitioned into 3 categories more or less equally distributed. The aim of this paper is to analyze the good (and bad) points of the competition, to fairly report the performances of solvers and emphasizes what we think was (or was not) a success, in terms of results and progress. The paper is organized as follows: in a first part, after recalling the main rules of the contest, we present the solvers, the benchmarks and the results of the contest identifying the winners<sup>1</sup>. Then, we analyse the first phase results and

---

<sup>1</sup> Due to the lack of space, we cannot report here exhaustive results. All results are exhaustively available on the competition web site

finally we discuss SAT2003 competition issues and propose some improvements for the next competition. We think that the competition is now mature enough to evolve in such way that it can fulfill its three foundations principles: motivate the field, identify new benchmarks and techniques and experimentally study solvers behavior.

### 1.1 Competition rules

The rules were available a few months before the competition started on the competition website. Briefly, the competitions runs in two phases. The idea is to use the first phase as a way to identify the best solvers and the hardest benchmarks for each categories, in order to use more time for the best solvers on the hardest benchmarks in a second phase. The first phase is thus important to provide exhaustive empirical results for analysis, while the second one is providing more limited results mainly dedicated to identify challenging benchmarks and award both solvers and benchmarks.

Like last year, we used the notion of series: benchmarks were partitionned into 3 categories (Industrial, HandMade, Random), then in series. The idea of a series was to group together benchmarks that are similar (same generator or same kind of problem, ...) to be able to identify solvers that can solve a wide variety of different benchmarks. A series was considered as solved if at least one benchmark in the series was solved. The solvers were ranked after the first phase according to the number of series and the total number of benchmarks solved. From those anonymized rankings, the competition judges, John Franco, Hans van Maaren and Toby Walsh (a new feature of the SAT2003 competition), decided which solvers should enter the second phase for each category. Note that only the best variant of the same solver was eligible for the second phase (which is also a new feature of the SAT2003 competition). After the second phase the winner was the solver that solved the greatest number of benchmarks. The competition awarded 6 solvers (3 categories of benchmarks in which we awarded the best complete solver (answers SAT+UNSAT) and the best solver that answered SAT (Complete + Incomplete solvers) and the smallest unsolved SAT (resp.UNSAT) benchmark. If a solver answered UNSAT on a proved (checked by a certificate) SAT instance, then that solver was declared “buggy”. During the conference, we presented buggy solvers as hors-concours solver and took only into account their SAT answer (because SAT answers are always checked with the certificate). For this report, we launched corrected versions of the buggy solvers so we report all their answers.

We also wanted this year to be as transparent as possible for submitters. During the whole process of the competition, solver submitters were able to check almost lively the progress of their own solvers. Results are available via the web, in a dynamic HTML and a raw format, which allows anyone to make its own analysis of the competition data.

From a practical viewpoint, the competition ran on two clusters of Linux boxes. The first one, from the “Laboratoire de Recherche en Informatique” (LRI, Orsay, France) was composed of 15 Athlon 1800+ with 1Gb of RAM. The second

one, from “Dipartimento di Informatica Sistemica e Telematica” (DIST Genoa, Italy), was composed of 8 PIV 2.4GHz with 1Gb of RAM. Both clusters were using RedHat Linux 7.2. During the first phase, industrial and handmade categories ran on LRI’s cluster while most of the random category ran in Italy. All the second phase was running on LRI’s cluster (while the QBF evaluation was running in Italy[14]).

## 2 The competitors

### 2.1 Solvers

Solver submission closed on February 14th but the competition really started one month later (in order for us to check input/output compliance and to receive corrected versions of the non compliant solvers). Late submissions were accepted *hors-concours*, if they complied with input/output requirements. Those solvers entered the first phase, which allowed them to compete with current state-of-the-art SAT solvers, without being awardable. The detailed list of the 34 solvers participating to the competition per category follows. Let’s begin with the description of the 28 complete non randomized solvers:

**berkmin** (release 561 and release 62) from Eugene Goldberg and Yakov Novikov [10]. Berkmin62 was awarded “best solver on satisfiable handmade benchmarks” at SAT’02 competition.

**bmsat** (BlueMoonSAT, hors concours) from Xiaowei Xu.

**compsat** (release 0.5) from Armin Biere. A bug was discovered during the contest.

We report here the fixed version, compsat-fixed, as a hors-concours solver.

**farseer** (release 1.00), from Martin Girard.

**forklift** from Eugene Goldberg and Yakov Novikov.

**funex** (release 0.2) from Armin Biere.

**jerusat** (release 1.1 va, vb, vc) from Nadel Alexander [2].

**jquest2** from Inês Lynce and João Marques Silva [18].

**kcdfs** from Gilles Dequen and Olivier Dubois [8, 7].

**limmat** (release 1.3) from Armin Biere. Limmat was awarded “Best solver on satisfiable industrial benchmarks” at SAT’02 competition.

**lsat** (release 1.1) from Richard Ostrowski, Bertrand Mazure and Lakhdar Sais [22].

lsat (v 1.0) participated to the SAT’02 competition.

**march sp/tt/xq** from Mark Dufour, Marijn Heule, Joris van Zwieten and Hans van Maaren. marchxq was found buggy. We report here the results of marchxq-fixed, the fixed version, as a hors-concours solver.

**oepir** Oepir (release 0.2) from Johan Alfredsson.

**oksolver** (hors concours) from Oliver Kullmann [13]. Oksolver was awarded both “best complete sat solver on random benchmarks” and “best solver on satisfiable random benchmarks” at SAT’02 competition. oksolver was not submitted to the competition but it is entering the competition as the winner of the previous competition.

We also report its performances on all categories.

**opensat** (release 0.44) from Gilles Audemard, Daniel Le Berre and Olivier Roussel [3].

**satnik** from Niklas Sörensson [9, 6].

**sato** (release 3.4) from Hantao Zhang [28].  
**satzilla** (release 0.9, v1 v2 v2-fixed) (v2s are hors concours) from Kevin Leyton-Brown, Eugene Nudelman, Galen Andrew, Carla Gomes, Jim McFadden, Bart Selman and Yoav Shoham [16, 15]. Authors asked us to consider a new version of satzilla2. They discovered, during the contest, that satzilla2 had bugs in its learning phase, that led to wrong choices of solvers during the contest. The new version is referred as satzilla2-fixed.  
**satzoo** (release 0.98 v0 v1) from Niklas Een [9].  
**tts** (ternary tree solver, release 1.0, hors concours) from Ivor Spence.  
**zchaff** from Yinlei Yu and Lintao Zhang [21, 29]. zchaff was awarded best complete solver for both industrial and handmade categories at SAT'02 competition.

The 2 complete randomized solvers were:

**siege** (hors concours) from Lawrence Ryan.  
**xqinting** (release 0.1) from Xiao Yu Li, Matthias F. Stallmann and Franc Brglez.

The 4 incomplete solvers were:

**amvo** from Sean Vogt and Andrew Machen.  
**qinting** (release 1.0) from Xiao Yu Li, Matthias F. Stallmann and Franc Brglez [17].  
**saturn** (release 2) from Steven Prestwich. [23].  
**unitwalk** (release 0.981) from Edward Hirsch and Arist Kojevnikov [11].

## 2.2 The benchmarks

We had 17 submitters, but a lot of benchmarks were submitted by the same people (BMC from IBM, generators for random/structured formulae). We also reused all SAT'02 challenges (except some from SAT'02 random challenges because there was too many of them). We tried to equilibrate the cpu effort on categories: Industrial with 323 benchmarks in 45 series, Handmade with 353 benchmarks in 34 series, Random with 317 benchmarks in 30 series. Like last year, we shuffled all the benchmarks to avoid syntactical recognition of benchmarks by solvers. Regarding random and handmade benchmarks, we only considered 3 benchmarks per point. That means that once we fixed all the parameters for the generator, we generated 3 (and only 3) benches with 3 different seeds. The new benchmarks were, for the industrial category (we have a total number of 323 benchmarks with, on average, 42703 variables, 163947 clauses and 452283 literals):

**hard\_eq\_check** submitted by E. Goldberg, 16 benchmarks, 1 serie.  
**addm** Submitted by J. Kukula, 6 benchmarks, 1 serie: equivalence check of two randomly structured adders with multiple addends.  
**li-exam, test** Submitted by R. Li, 12 benchmarks, 2 series: formulated from unit covering problem of logic minimization of logic ircuit benchmark test4  
**ferry, gripper, hanoi** Submitted by Maris, 24 benchmarks, 3 series: based on TSP (tunable satplan). A DIMACS generator for planing problems in PDDL format  
**l2s** Submitted by V. Schuppan 11 benchmarks, 1 serie: Bounded model checking of Bwolen Yangs collection of benchmarks for SMV in conjunction with a method to convert liveness to safety.

**zarpas** Submitted by E. Zarpas: a lot of (large) BMC formulas. We selected only a subset of 223 benchmarks in 25 series.

For the handmade category (we have a total number of 353 benchmarks having, on average 3233 variables, 37355 clauses and 154485 literals):

**sgi** Submitted by C. Anton: Subgraph isomorphisms, 120, benchmarks, 8 series.

**graphs** Submitted by R. Bevan: Formulas based upon various graphs as first described by Tseitin in 1968, 56 benchmarks, 7 series. Graphs used include margulis graphs, modified margulis graphs and various other well known graphs (i.e. hypercube).

**quasigroups with holes** Submitted by C. Gomes, 19 benches, 2 series.

**chesscolor** Submitted by P. Lardeux, 5 benchmarks, 1 series: These benchmarks encode an original chessboard coloring problem. A  $n^2$  chessboard must be colored with  $k$  colors such that the four corners of every rectangles included in the chessboard are not all of the same color.

**multLin, mm** Submitted by K. Markstrom, 17 benchmarks, 1 series: CNFs coding various instances of fast matrix multiplication, both for square and non-square matrices. The encoding of the problems is intended to be efficient.

**genurq** Submitted by R. Ostrowski, 10 benchmarks, 1 series: Modified genurq problem, all satisfiable.

**clusgen** Submitted by A. Slater 9 benchmarks, 3 series: Randomly generated clustered 3SAT problems.

**hwb** Submitted by T. Stanion, 21 benchmarks, 1 series: The generator build two different circuits computing the hidden weighted bit function and combines them into an equivalence checking problem.

We mainly restricted the random category benchmarks that looked like uniform random benchmarks. We have 317 benchmarks with on average 442 variables, 4164 clauses and 25401 literals:

**balanced, hidden** Submitted by C. Moore, 111 benchmarks, 8 series: formulas with 1 or 2 hidden assignments and balanced literals: satisfiable but hard.

**hgen\*** Submitted by E. Hirsch, 57 benchmarks, 4 series: generators based on hgen2 (last year winner)

**gencnf** Submitted by G. Dequen, 63 benchmarks, 7 series: k-cnf uniform generator at threshold for  $k = 4-10$ .

**uniform** Submitted by L. Simon, 45 benchmarks, 5 series, a uniform generator for  $k=3$ , one serie per constant ratio (only one serie at threshold).

### 3 First phase results

The first phase began on March 19th and ended on April 14th. One of the new features of the SAT'03 competition was the ability for each competitor to follow his solver via a web page that was updated daily. All the traces were available, and suspicious launches (due to hardware/network error, ...), detected by us or by the competitors, were ran again. The two solvers compsat and marchxq were found incorrect ("buggy") during the first phase due to minor bugs. We present here the result of their fixed version, compsat-fixed and marchxq-fixed. satzilla2-fixed is a new version of satzilla2 with improved learning scheme.

In this section, we discuss the results of the solvers in each category. We will not discuss here the results of the SAT subcategory in both the industrial and handmade categories, since there are not enough benchmarks to draw any conclusions. Those results are available on the competition website.

### 3.1 First phase on Industrials

One can note that there was many zchaff-like solvers this year (13 solvers, 17 variants, against 4 solvers last year). This can be explained by the huge interest in the practical resolution of SAT in the EDA community since many of the new solvers are coming from that community, and by the fact that zchaff was a breakthrough in the community in the last two years. For instance, 3 out of 4 of the solvers awarded last year were similar to zchaff (including the original).

The figure 1 illustrates the results of the complete solvers on all the benchmarks of the industrial category. This kind of representation, used in the previous competition report [26], allows to check the choice of the cpu timeout and to have clues about solvers behaviors on all the benchmarks in a given category. Last year winners in the category were zchaff and limmat. Those solvers end a group of zchaff-like solvers led by forklift. If we take those solvers as the references for state-of-the-art SAT solver in 2002 then it seems clear that some progress was made in the industrial category (half of the submitted solvers are more efficient than these progress witnesses!).

But the picture is not so clear. Let us take a look at berkmin62 curve. One can note that not that many solvers are stronger than berkmin62. Let us recall that the version of berkmin62 that entered the SAT'02 competition had a bug that prevented it to solve hundreds of benchmarks[26]. The version that entered the SAT'03 was fixed. So it is reasonable to consider berkmin62 as a progress witness for the year 2002. Thus, the progress made is not that important: three solvers only performed better than berkmin62 and two of them are from the same authors, one solver being a variant of one of the engine of berkmin62, while the other one is a new solver that can be seen as an extension of berkmin62 with binary clause reasoning.

So, the good point is that a lot of new solvers were proposed and showed relatively good performances for their first public evaluation. But the pernicious effect of the competition is that, in the run for good performances, many solvers are very similar (zchaff-like).

### 3.2 First phase on Handmade

The figure 2 illustrates the results of the complete solvers on handmade benchmarks. One solver, lsat, clearly shows a particular behavior. This solver doesn't implement only resolution, but can handle (even long) chains of equivalency clauses and can more generally retrieve from a CNF some boolean function on which it applies some simplifications [22]. That point gives it the ability to quickly solve tricky formulas, commonly used in this category (where formulas are build to defeat resolution-based algorithm). Beside this very particular

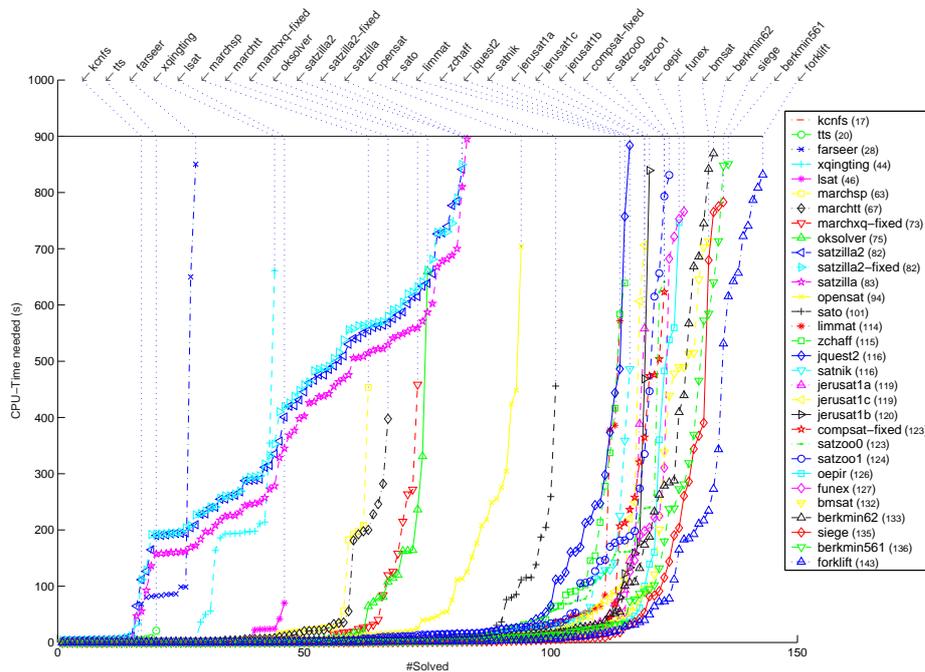


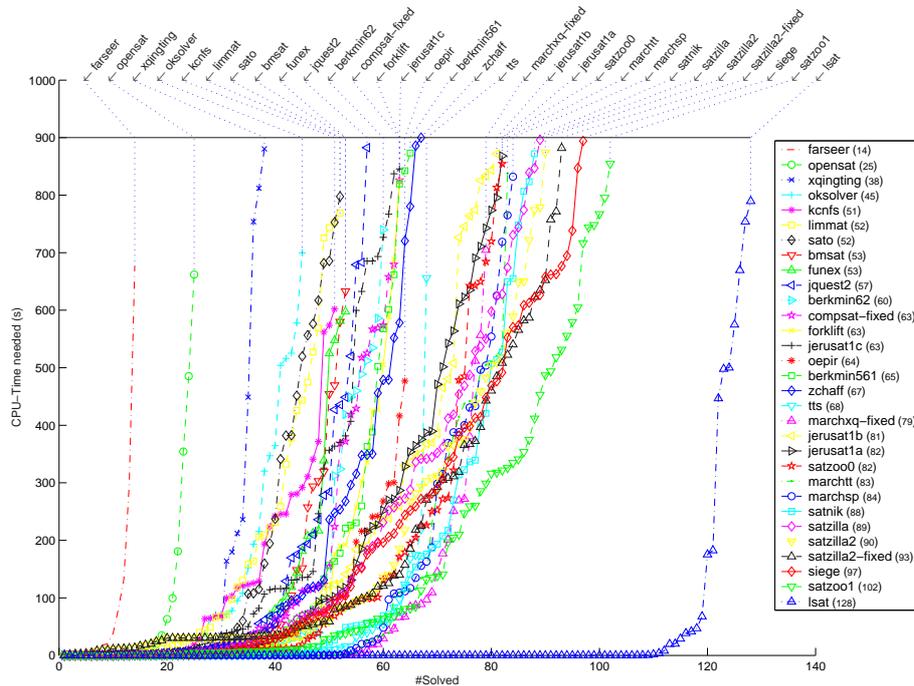
Fig. 1. Number of instances solved vs. CPU time for complete solvers on all industrial benchmarks

solver, the figure shows relatively efficient solvers between marchxq-fixed and satzoo1 on this category.

### 3.3 First phase on Random

Results for the two random categories are given on figures 3 and 4. On the first figure, one can easily identify all complete solvers good for random instances (between oksolver and satzilla2-fixed). The curves have a relatively slow growth (certainly due to the fact that random instances scales smoothly), especially for knfs, which seems to scale well. This is not the same picture for satzilla, which present the same kind of curve than lsat on HandMade: a lot of benchmarks are solved quickly, and then the curve have an important growth. On random SAT instances, unitwalk and satzilla2 present the same kind of curves.

**The case of satzilla** The solver takes an unusual, “portfolio” approach to the SAT problem. It uses machine learning models to select among a set of existing SAT algorithms, and then runs the algorithm predicted to do best. For satzilla 0.9, these include: 2clseq, limmat, jerusat, oksolver, relsat, sato, satzrand, and zchaff. It then runs the algorithms with the lowest predicted run time. The techniques for predicting run time of algorithms are described in [16].

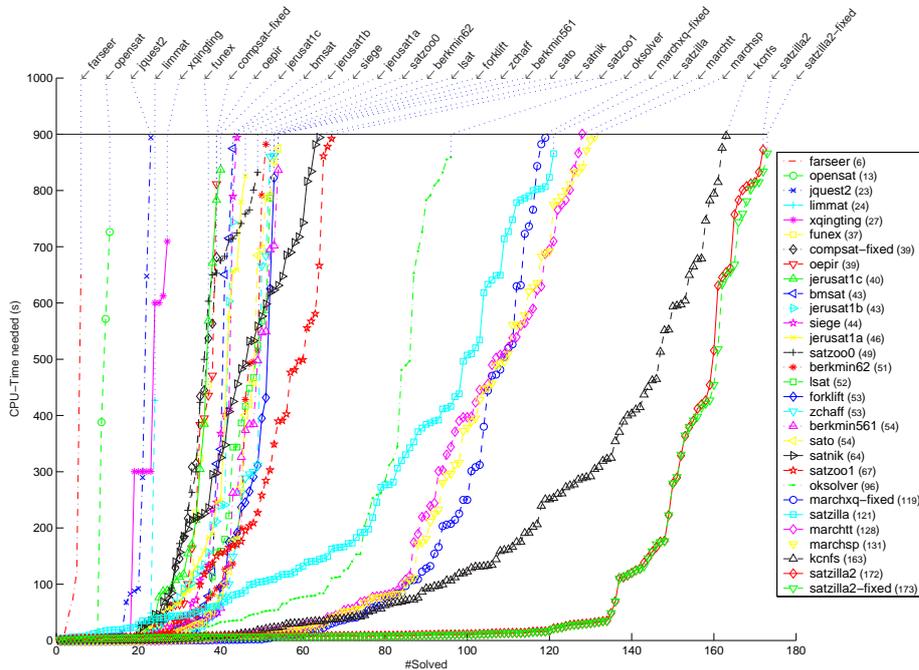


**Fig. 2.** Number of instances solved vs. CPU time for complete solvers on all handmade benchmarks

The portfolio approach (with experimental results examining a different problem domain) is elaborated in [15]. Hors Concours solver satzilla2 (submitted after the deadline, on February 28) uses two additional complete solvers, eqsatz and heerhugo. It also executes the autoSat algorithm for a short time before starting any other computation, and is thus able to filter out easy satisfiable problems (this point can explain the kind of curves we observed).

Because this approach is unusual and represent one of the main new approach in this year contest, we conducted a few analysis on satzilla2-fixed to understand how it behaves. We provide for each solver composing satzilla2 how many times it was used and how many times it solved the benchmarks (number in parenthesis). The missing solved benchmarks are due to AutoSAT.

- On industrial instances, it launched jerusat 13 (3) times, zchaff 64 (49) times, eqsatz 1 (0) time, 2clseq 27 (4) times, satzrand 4 (0) times and limmat 2 (0) times.
- For HandMade, the choices were jerusat 93 (42), relsat 1 (0), eqsatz 26 (6), 2clseq 5 (0), satzrand 39 (9), limmat 62 (1), oksolver 26 (1), heerhugo 7 (0) and sato 10 (3).
- On random instances, the choices were jerusat 1 (1), zchaff 12 (2), eqsatz 18 (18), 2clseq 13 (1), satzrand 124 (37) and oksolver 35 (1). The choices were



**Fig. 3.** Number of instances solved vs. CPU time for complete solvers on all random benchmarks

different for each category and one may notice that zchaff was chosen 12 times on random instances (and even gave 2 results).

Note that eqsatz was launched 18 times on random problems and succeeded each time. This is because the hardnm benchmarks from C. Moore contains equivalency chains that can be solved in linear time by specific solvers such as eqsatz or lsat. For such a reason, the behavior of satzilla2 on those “random” benchmarks must be taken with care.

### 3.4 Solvers entering the second phase

This year, the three judges decided, from an anonymized version of the rankings presented table 1 which solvers should enter the second phase of the competition. Table 1 also emphasizes the solvers that were alone to solve at least one benchmark (so called State-Of-The-Art Contributor, SOTAC, in [27]) during the first phase. In addition to the one reported in the tables, marchxq-fixed (Industrial) and compsat-fixed (Industrial, HandMade) are also SOTAC.

### 3.5 The hardest benchmarks

The pool of benchmark remained unsolved during the first phase is available on the web. We encourage developers to focus on this set of hard bench-

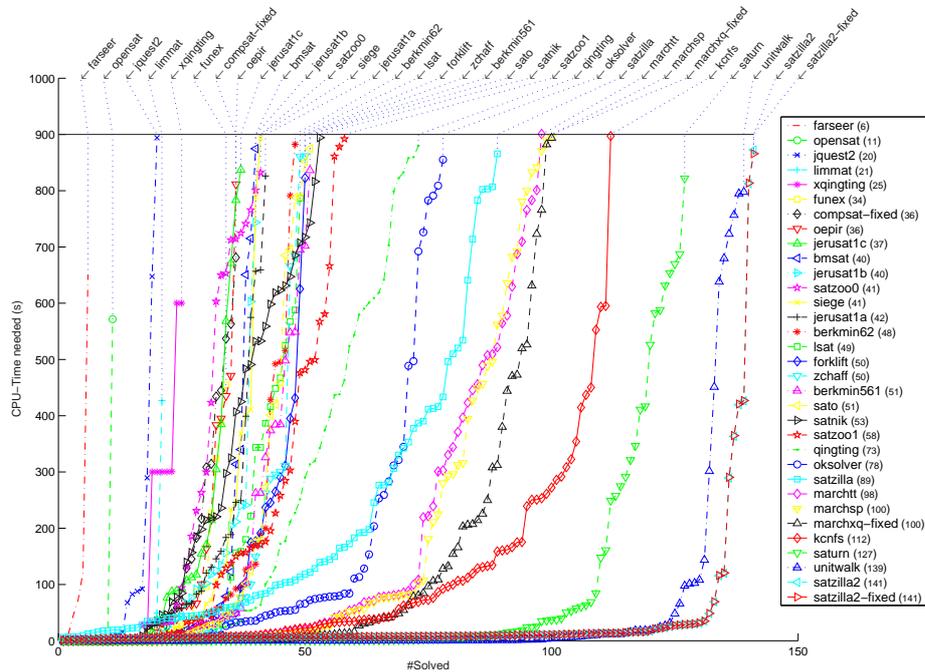


Fig. 4. Number of instances solved vs. CPU time for all solvers on SAT random benchmarks

marks. We have 174 benchmarks (from 16 series) in the Industrial category (mean (min-max) sizes: #clauses=222811 (7814-3113540), #variables=61322 (2269-985042), #literals=627908 (21206-7976612)), 165 benchmarks (from 18 series) for the HandMade one (#clauses=54422 (832-630000), #variables=4518 (75-270000), #literals=263992 (2818-1350000)) and 91 (from 17 series) for the Random category (sizes: #clauses=7838 (391-36000), #variables=286 (38-700), #literals=61515 (888-360000)). Among the unsolved benchmarks, we often find large industrial formulas and, in the random category, k-cnf with large k.

## 4 Second phase

We selected 269 benchmarks that went to second phase. To give an idea of their sizes (you can compare numbers with the one in the previous section 3.5), we selected 72 benchmarks (in all the 16 series) in the Industrial category (sizes: #clauses=324353 (7814-3113540), #variables=92473 (2269-985042), #literals=932932 (21206-7976612)), 137 (in all the 18 series) in the HandMade category (sizes: #clauses=52561 (832-630000), #variables=6027 (75-270000), #literals=238815 (2818-1350000)) and 60 (in all the 17 series) in the Random category (size : #clauses=6258 (391-30960), #variables=329 (38-700), #literals=46908 (888-309600)).

Solvers	#series	#benchs (comments)	Solvers	#series	#benchs (comments)
<b>Industrial</b>					
Complete solvers on all benches			All solvers on SAT benches		
forklift	34	143 (SOTAC)	forklift	9	33
berkmin561	32	136	berkmin62	8	32
satzo01	31	124	zchaff	8	31
oepir	30	126	oepir	8	31
funex	29	127 (SOTAC)	jerusat1b	8	31
satnik	29	116	funex	8	31
jquest2	29	116	satnik	8	30
zchaff	29	115 (SAT'02 winner)	limmat	8	30 (SAT'02 winner)
jerusat1b	28	120	satzo01	7	30
limmat	28	114	jquest2	7	29
<b>HandMade</b>					
Complete solvers on all benches			All solvers on SAT benches		
jerusat1b	22	81	satzo01	11	37 (SOTAC)
satzo01	21	102	jerusat1b	11	35
satzilla	21	89	satzilla	11	25
satnik	20	88	forklift	10	30 (SOTAC)
marchsp	19	84	oepir	9	31 (SOTAC)
lsat	15	128 (SOTAC)	satnik	9	30
			berkmin561	9	30 (v62 was SAT'02 winner)
			saturn	7	29 (SOTAC)
<b>Random</b>					
Complete solvers on all benches			All solvers on SAT benches		
kcdfs	25	163 (SOTAC)	kcdfs	18	112
satzilla	19	121	unitwalk	15	139 (SOTAC)
marchsp	18	131	saturn	14	127
oksolver	16	96	satzilla	14	89
satnik	13	64	marchsp	13	100

Table 1. Ranking of solvers entering the second phase.

#### 4.1 Winners: the solvers!

The table 2 summarizes the whole second phase of the contest. Here, the cpu timeout was set to 2 hours (instead of 15mn). The 4 winners are in bold face in each category. Briefly, **forklift** won both (SAT+UNSAT and SAT) Industrial categories, **satzo01** won both HandMade categories. **kcdfs** won the Random SAT+UNSAT category and **unitwalk** the Random SAT category.

If this is not a surprise for forklift, according to the results of the first phase. Now, the fact that kcdfs won a price may be a good indicator that our partitioning may be not so bad. Simply because kcdfs was devoted to solve such instances. We also want to point out that this is the first time an incomplete solver (unitwalk) win the competition. Last year, oksolver won both the Random categories prizes. oksolver was even not qualified for the second phase on Random SAT instances (it was on SAT+UNSAT).

At last, one have to point out that winners are often determined by very small difference. This is not really surprising since the shapes of curves are almost vertical on figures 1–4 for these solvers, but may be disappointing since after one month of contest, the winner was determined on its performances on a very few benchmarks.

Solver	#benchs	Solver	#benchs	Solver	#benchs
Industrial Complete		Industrial SAT		HandMade Complete	
<b>forklift</b>	<b>12</b>	<b>forklift</b>	<b>4</b>	<b>satzo01</b>	<b>9</b>
berkmin561	11	jerusat1b	3	lsat	7
satzo01	5	satnik	3	satzilla	6
jerusat1b	5	funex	1	satnik	3
satnik	4	zchaff	1	jerusat1b	1
zchaff	4	oepir	1	marchsp	0
funex	3	berkmin62	1		
oepir	1	satzo01	1		
jquest2	0	limmat	0		
limmat	0	jquest2	0		
Random Complete		Random SAT		HandMade SAT	
<b>kcnfs</b>	<b>12</b>	<b>unitwalk</b>	<b>12</b>	<b>satzo01</b>	<b>5</b>
satzilla	6	kcnfs	4	satzilla	3
oksolver	4	saturn	2	forklift	3
marchsp	3	satzilla	2	berkmin561	1
satnik	0	marchsp	0	saturn	0
				oepir	0
				satnik	0
				jerusat1b	0

**Table 2.** Ranking of solvers after the second phase

## 4.2 Winners: the benchmarks

In each categories, we selected the smallest benchmarks to use during the second phase, according to their number of literals, in order to award the smallest unsolved benchmarks in both SAT and UNSAT categories.

The winner for the smallest unsat benchmark was `random/hirsch/hgen8/-hgen8-n260-01-S1597732451`, a benchmark submitted by E. Hirsch. It was renamed `sat03-885` during the contest. This unsatisfiable benchmark contains 391 clauses, 260 variables and only 888 literals. In comparison to last year results, the smallest UNSAT instance moved from the HandMade to the Random category, which can be explained by two reasons. First, `lsat` has dropped away, during the first stage, all small benchmarks based on equivalency clauses and second, the category of this benchmarks may not be the good one.

On the SAT category, we didn't award any benchmark. The winner was `random/simon/sat02-random/hgen2-v400-s161064952` (renamed `sat03-1681` this year), a benchmark submitted by E. Hirsch last year. It contains 1400 clauses, 400 variables and 4200 literals. That benchmark was solved in the second phase last year, but not this year. Such a phenomenon may be due to bad luck this year for randomized instances, or simply to the re-shuffling of the instance. We'll discuss this kind of problem in the section 5.2.

## 5 Discussion

### 5.1 Other rankings

To check the strongness of the competition final results, one may wants to try other methods to rank solvers and to compare the rankings. This is also a way to analyse the results of the first phase. In the contest, we tried to forget about

Solver	Total Time (s)	#solved benches	Solver	Total Time (s)	#solved benches
Complete on Industrials			Complete on HandMade		
forklift	171240	143	lsat	207371	128
siege	175168	135	satzo01	241664	102
berkmin561	176088	136	satnik	249446	88
bmsat	178445	132	siege	250043	97
berkmin62	179121	133	marchtt	250711	83
Complete on Random			All On SAT Random		
satzilla2	147208	172	satzilla2	35402	141
kcdfs	165655	163	unitwalk	38572	139
unitwalk	166372	139	saturn	52136	127
saturn	179936	127	kcdfs	67493	112
marchsp	192688	131	marchsp	83030	100

**Table 3.** Top-5 ranking of solvers *à la satex*

the cputime, by counting only the number of series and benchmarks solved. Of course the cpu time had a direct impact on the competition due to the timeout, but it is important to try to characterize efficient solvers, in term of cpu time.

**Ranking *à la satex*** One way to rank solvers is to use cumulative cpu time, as it is done in SatEx [25]: just count the total cpu time needed for one solver to solve the whole category of benchmarks, using the timeout as the default penalty. Note that the timeout was 15 mn, which is not a big penalty, so a fast solver can be ranked above a solver that solved more benchmarks.

Table 3 summarizes the top-5 solvers in each category that we studied in this paper. Forklift is still a clear winner (siege was hors-concours) on Industrial benchmarks. On HandMade benchmarks, such a ranking would have awarded lsat instead of satzo01. Lsat was able to quickly solve a large amount of benchmarks without search, but when the DLL engine of lsat was needed (on hard instances like in the second phase), lsat was not able to beat satzo01. On random benchmarks, satzilla2 (which is hors-concours) would be ranked first, which is in part due to the 18 hardnm benchmarks outlined before.

**Relative efficiency of solvers** One of the hard thing to handle for ranking solvers is that only partial information is available. We have to use a cputime timeout and, what ever value we choose, we'll never have a complete picture of all solvers cpu time on all benchmarks. One possible way is to compare only pairs of solvers  $(X, Y)$  on the subset of benchmarks they both solves (if  $s(X)$  is the set of benchmarks solved by  $X$ , then we compare  $X$  and  $Y$  on their respective performances on the set  $s(X) \cap s(Y)$ ). When doing this, we have a strong way of comparing the relative efficiency (RE) of  $X$  and  $Y$  :  $re(X, Y) = s(X) \cap s(Y) / s(Y)$  gives the percentage of instances of  $Y$  that  $X$  solves too. Let us write now  $cpu(X, b)$  the cpu time needed for  $X$  to solve all the benchmarks in  $b$ , without any timeout. Because there was a timeout in the competition, only  $cpu(X, s')$ , with  $s' \subseteq s(X)$  are defined here for the solver  $X$ . Now, we can compare the relative efficiency of  $X$  with  $Y$  by computing  $crr(X, Y) = cpu(X, s(X) \cap s(Y)) / cpu(Y, s(X) \cap s(Y))$ . This last measure is called here the

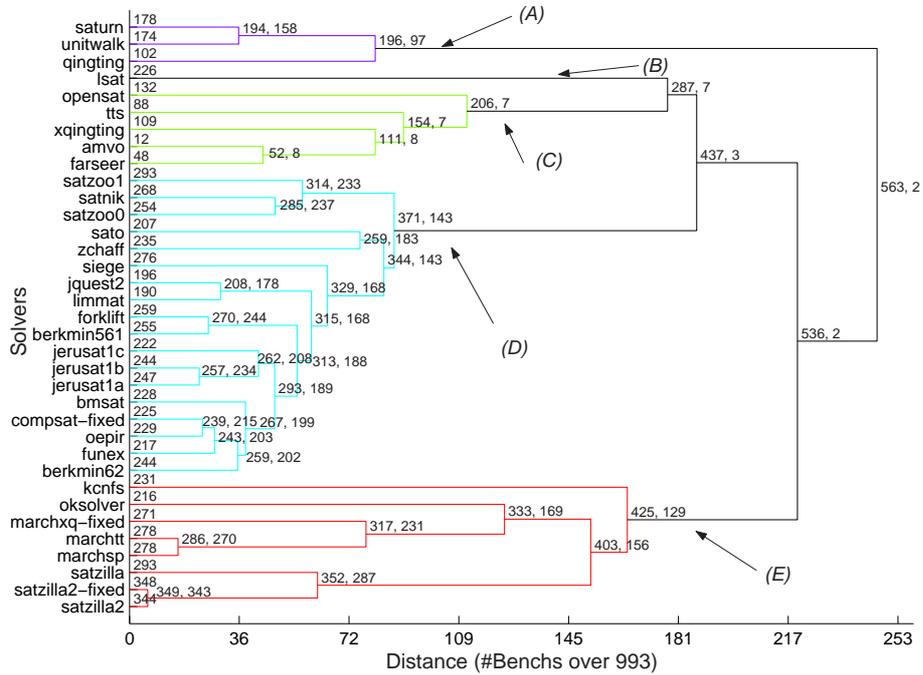
cputime reduction ratio (CRR), and means that, on their common subset of solved benchmarks,  $X$  needs only  $crr(X, Y)$  percent of the time needed by  $Y$ . Now, to summarize all the possible values, we average these two measures over all the possible values for  $Y$ , while keeping  $X$  fixed.

In the Industrial (complete) category, we have the following values for  $re$  : forklift (98.9%), berkmin561 (98.7%), siege (98.6%), berkmin62 (98.4%), oepir (97.6%), funex (97.4%), bmsat (97.2%), compsat-fixed (96.8%), satzoo0 (96.4%) and satzoo1 (96.3%). The  $crr$  values are forklift (23%), siege (28%), bermin561 (46%). funex and oepir needs only 55% of the time of the other solvers. The  $crr$  values of solvers then jumps to more than 70% for the other solvers. It is interesting to remark that both forklift and siege are very efficient solvers ( $crr < 25\%$  and an impressive value for  $re$ ). Satzoo0 is also interesting because it has a high  $re$  and a  $crr$  at 164%. This is the same kind of remark for jquest2, which has a  $re$  of 94.5% and a  $crr$  at only 388% (certainly due to the java penalty). In a sense, we think that such a measure may emphasize a robust solver (if  $s(X)$  is large enough): zchaff has a  $re$  of 91.1% and a  $crr$  of 87%, and limmat a  $re$  of 93.9% and a  $crr$  of 138%. Thus, jquest2 and limmat are slower than zchaff, but can on average solve more benchmarks.

In the HandMade (complete) category, the  $re$  values are satzoo1 (88.9%), siege (88.6%), satnik (86.3%), satzilla2-fixed (86.2%), satzoo0 (84.9%), satzilla2 (84.9%), jerusat1a (83.8%), jerusat1b (83.2%) and satzilla (83.1%). It then drops to less than 73% for the other solvers. The best  $crr$  are for lsat (56%), satzoo1 (83%), marchxq-fixed (90%) and forklift (97%). All other solvers have a ratio greater than 100%. For information, lsat has a  $re$  of 72.6%. This is a low value, but this kind of measure can be misleading: in this case, lsat did not solve benchmarks solved by other solvers, which is reflected by the small number of series solved, while solved quickly others (which is reflected by the total number of benchmarks solved, being SOTAC, etc).

In the Random (Complete) category, we have the following  $re$ : satzilla2-fixed (95.4%), satzilla2 (95.3%), marchxq-fixed (90.4%), marchsp (83.3%), marchtt (84.0%) and kcnfs (83.9%). The percentage then jumps to less than 75%. The best  $crr$  are satzilla2-fixed (only 26%), satzilla2 (27%), kcnfs (39%), marchxq-fixed (47%). The ratio then falls to more than 88% (with oksolver).

At last, in the Random (SAT) category, we have the following  $re$ : satzilla2 (92.8%), satzilla2-fixed (92.8%), marchxq-fixed (87.0%), marchsp (82.2%), marchtt (80.1%) and kcnfs (78.9%). The best  $crr$  are different here: unitwalk, have only a  $re$  of 71.4% but a very remarkable  $crr$  of 6%! saturn has a  $re$  of 71.1% and a  $crr$  of 28%. That means that unitwalk and saturn solve problems different from other solvers, and solve the common problems very quickly. In the order, after these two incomplete solvers, the  $crr$  are sazilla2-fixed and satzilla2 (29%) and kcnfs (81%). The  $crr$  then falls to more than 100%. In some sense, unitwalk present the same kind of results than lsat: a low  $re$  value that indicates here an original method. If  $crr < 1$  then this method is fast, and, at last, the method is *interesting* if  $s(x)$  is large enough. This is clearly the case for lsat and unitwalk.



**Fig. 5.** Clusters of all solvers on all benchmarks. Solver that closely solves sets of benchmarks are close together in the cluster. The height of nodes in the tree indicates the average distance of the two clusters at each considered branch of the tree. Number at the left indicates the number of solved instances for each solver, and the couple of number at each cluster link indicates respectively the number of common benchmarks solved by at least one member of the cluster and the number of benchmarks solved by all the members of the cluster. 5 main cluster A–E can be easily identified

**Clustering of solvers according to their performances** We automatically clustered the solver according to the set of common benchmarks they solved. The clusters are represented with a tree in figure 5.

Related solvers all belong to the same cluster. For instance, different versions of the same solver are very close in the cluster (see marchXX, jerusat1X and satzillaX). The cluster (A) denotes incomplete solvers. The cluster (B) contains only one solver: the singular lsat. The cluster (C) is made of medium strength complete solvers (the number of common solved benchmark is very low due to amvo’ performances). Then we have this large cluster (D) of zchaff-like solvers. This cluster contains solvers that are good on industrial benchmarks. One may remark here that 143 benchmarks are solved by all the solvers of (D), while 371 benchmarks (over 993) are solved by at least one of these solvers. All the solvers in this cluster have very close performances. In the last cluster, noted (E), one can identify the complete solvers strong on random instances.

## 5.2 Randomized or not: the *lisa* syndrom

It was brought to the fore in [17] that solvers may have drastically different behaviors on a range of problems generated from one original benchmark by randomly shuffling both its clauses and renaming its literals. The shuffler used for the competition generates members of the PC-class of [17]. However, this very important addressed in the competition. This year, the same benchmark (a SAT'02 challenge called *lisa*) was placed in two different categories by error. Before we detected the problem, it was shuffled twice (one for *lisa* in Industrial and one for *lisa* in HandMade with different seeds) and solvers ran on both instances. Here are some examples of solvers that showed different behaviors on those two benchmarks (“-” stands for *unsolved*): *jerusat1a*: (-, 40.71) ; *jerusat1b* (236.85, -) ; *oepir* (-, 12.62) ; *siege* (67.43, -), *satzo01* (-, 287.12) and *satzilla* (196.8, 727.52).

Such a result can question the validity of the competition. However, such a thorny problem may be smoothed by the large number of benchmarks used and we strongly trust the main picture provided by the rankings we reported for the first phase. This is not the case for the second phase. Considering the small number of benchmarks solved, being lucky can make the difference.

## 5.3 Progress or not?

One of the most important question when assessing the results of the competition was whether or not substantial progress was made. Our first thought was “yes”, since several benchmarks unsolved during 6 hours last year were solved within 15 mn this year.

Let's take a closer look at each category. Progress on the Industrial category was already discussed in section 3.1. In the handmade category, *lsat* shown very good performances on many structured benchmarks. However, *lsat* participated last year and became hors concours because of a bug. This year version is mainly last year version without the bug. So it is hard to see here any real progress since last year.

In the random category, the picture is a bit different. Some new solvers were submitted this year. Last year winner, *oksolver*, was outperformed in the satisfiable category. The most interesting thing concerning that category is the presence of the hors concours *satzilla2*, a portfolio solver, that compared favorably against the awarded complete solver and did slightly better than the incomplete solver awarded for satisfiable benchmarks.

## 5.4 Benchmark categories: back to the roots

Some improvements were made concerning the way solvers were evaluated this year: only the best variant of a given solver was eligible for the second phase, solver submitters were able to follow almost in real time their solver during the first phase of the competition, the judges chose the solvers that should enter the second phase, which was more flexible and fair than using an arbitrary rule

(top 5 solvers...), etc. But nothing really changed regarding the way we handled benchmarks. We need to focus on this issue for the next year competition.

The idea of creating three categories of benchmarks, namely industrial, hand-made and random, was to cover the different interests of researchers in the practical resolution of SAT. The industrial category was aimed to address the practical resolution of SAT problems for SAT-based reasoning engines. Ideally, that category should reflect the strengths or weaknesses of the solvers as a SAT component, in a “real/industrial” setting. Bounded model checking[5, 4] or more generally all the applications from the Electronic Design Automation community[19], planning[12], cryptography[20] were good candidates for such category. Up to now, we received mainly benchmarks from the EDA community. Those benchmarks really reflected the kind of problems to be solved by an embedded SAT engine. Unfortunately, since we shuffled the benchmarks, we did not really evaluate the solvers in “real conditions”. The order of the clauses, the identifiers of the variables may provide some structural information about the initial problem and some solvers built to solve those kind of problems were designed to handle that information. By shuffling the benchmarks, we hide that information. Thus a solver can have a very different behaviour during the competition and “in situ”, which is in contradiction with one of the aim of the contest. We have three solutions. First, we keep all things like this and we let submitters explicitly write heuristics in their solver that take shuffling into account. If this allows to emphasize the existence of implicit heuristics based on the syntactical order of the formula (which is a good point), this kind of solution is weird in some sense: solvers will spend time trying to rearrange the formula! We can also avoid shuffling for this category, but then solvers may be too much tuned for benchmarks, casting doubts on the results on challenging, publicly available, benchmarks. The last choice may be to launch the solvers on the original benchmark and X shuffled benchmarks of the same benchmarks (this solution may also address the lisa syndrom).

The random category was designed to evaluate the progress made on a class of theoretically well known and hard (NP-complete) problems: random k-sat. There are not that many ways to generate random k-SAT problems, so most of the benchmarks used for that category were based on some uniform k-SAT generators. However, we received some generators able to build problems looking syntactically like random uniform problems, but are forced SAT or UNSAT (like all the generators from Edward Hirsch). The question is: were those generators in the right category? Even worst is the case of the hardnm benchmarks related previously. Theoretical and practical work on random k-SAT is based on uniform k-sat. Furthermore, some problems such as quasigroup with holes, or clusters of random k-SAT formulas were in the handmade category. So the random category would better be renamed “uniform random” category and all the non uniform random generators should move to another category. Whether that category should be a new one or should be the handmade category is not yet decided.

Industrial			Handmade			Random		
solver	#series	#benchs	solver	#series	#benchs	solver	#series	#benchs
forklift	30	110	satzilla	13	64	kcnfs	9	51
berkmin561	29	105	siege	13	60	satzilla2-fixed	6	32
siege	29	104	satnik	13	58	satzilla	6	32
berkmin62	29	101	jerusat1a	13	49	marchsp	6	31
bmsat	28	101	jerusat1b	13	46	marchtt	6	30
			satzoo1	12	65	marchxq-fixed	4	19
			lsat	11	107			

**Table 4.** Top ranking of solvers on hypothetical UNSAT categories. Please do notice that hors-concours solvers appear here.

The handmade (crafted) category was aimed to point out new techniques not based on resolution. Finding an efficient technique to tackle SAT problems not based on resolution is one of the ten propositional challenges proposed in [24]. That category was also a reminder that if some solvers can solve huge benchmarks in the industrial category, some small benchmarks are still out of reach of those solvers. The solvers awarded so far in that category (zChaff, Berkmin and satzoo) are using resolution. So the category is awarding very efficient solvers (typically with the greatest number of decisions per second) rather than clever ones. Note that this year, an hybrid solver recovering potential structural knowledge and using it to simplify the initial problem solved the greatest number of problems in that category, but was not awarded. Awarding the most efficient solver is of interest, pointing out “clever” solvers would be a plus for the competition.

At last, we observed in both industrial and handmade category very few satisfiable benchmarks. One way to solve this problem is to filter satisfiable benchmarks from a pool of benchmarks submitted to the competition. In that case the benchmarks are needed long before the solvers, which means to change the way we benchmarks are submitted. Another solution is to ask benchmarks submitters a balanced number of SAT and UNSAT benchmarks.

**An UNSAT category?** It is clear that incomplete (1-sided SAT) solvers can only be compared on the basis of SAT instances. A question that could be raised is that, because we have the complete (SAT+UNSAT) and the SAT subcategories, this choice may disadvantages solver that are only good on UNSAT instances. Furthermore, it is possible that, in the next few years, 1-sided-UNSAT solvers (that can only solve UNSAT instances) may arise. For both reasons, we decided to simulate the result of the first phase on a new UNSAT subcategory. The results are reported on table 4. They are similar to the results of SAT+UNSAT categories on Industrial and Handmade instances. This may be due to the little number of SAT instances in these categories. However, on Random instance, we can see that the marchXX family of solvers are very close now to the satzillaX family. This is certainly due to the fact that satzilla performs a local search of the solution and quickly solves easy SAT instances. This advantage over the marchXX solver is now over and seems to prove the relative

efficiency of the marchXX solvers on UNSAT Random instances, which was not clear in the SAT+UNSAT results.

## 6 Summary of proposals for SAT'04 competitions

Despite the huge amount of data collected during the contest, it is impossible to really understand why a solver is better than another. In order to empirically study solvers (e.g. find hypothesis on their behaviour, ...), we will probably ask submitters next year to output statistics about their runs. XML output format could be great for automated analysis of results (each submitter may then provide its own ontology) but we are thinking of something even simpler. Such information may lead to identify the solver that explore the smallest search space, the fastest solver on unit-propagation, ... and then to begin to really understand experimentally what's going on.

Now, one of the problem this year was to choose the right category for all benchmarks. A lot of different solutions exist. First, one could simply use the huge amount of data in the first phase to automatically (*e.g.* statistically) clusterizes the benchmarks into different categories. This is for instance done by the authors of *satzilla* that compute 60 polynomial time features per instance to decide which solver to use on which benchmark. However, we may have a large number of different categories and we think that is very usefull to keep the process in the hands of human. An automatic partitionning would imply that a given bench may change of category from one year to another one. Moreover, the original notion of categories is somewhat well understood and corresponds to some expert knowledge about benchmarks. Another solution would be to ask for the help of judges to classify benchmarks, according to some description given by the author.

How to “officially” point out new (promising) methods? The quest for fast-solvers is the pernicious effects of the contest and we have seen that this year, not so much new methods have been proposed. Because one of the aim of the contest was to promote new clever methods, it could be nice to have some kind of special prizes, given by the judges to any solver, based on the analysis of the first phase results. Of course, this solver may not be the fastest and the more efficient one, but should be the most “promising” one.

The second phase is also a point which may be questionned. We have seen that winners may simply be lucky on a very few instances. In some sense, it is frustrating to spend years of cpu time and to determine the winner on its performances on less than 5 instances (where luck can take a large part). Moreover, all the important data are gathered during the first phase. So, could we reconsider the second phase? We are thinking of two alternative solutions: determine the winner according to its performances on the first phase or using some kind of tournament in the second phase, where each step would ensure a clear winner between 2 solvers.

Briefly, a lot of other points may be discussed before next year competition:

- To begin the competition earlier, just discard solvers not complying with the input/output requirements.
- Do not run 3 time randomized solvers since deterministic solvers can be lucky too (cf the lisa syndrom).
- Collect benchmarks in a different way: we need more control over benchmark submission in order to ensure a meaningful competition.
- Allow benchmarks submitters to follow “live” the solvers on their benchmarks.

The most important issue for next the next competition will be the notion of classes of benchmarks instead of just one instance, to obtain strong result on solvers performances. A reasonable solution is to run the solvers on the original benchmarks (without shuffling them) and on a set of  $X$  shuffled version of that benchmarks,  $X$  depending on the number of benchmarks, solvers and computers available. Such approach provides a view of the solvers “in situ” and assesses their robustness.

## 7 Conclusion

The competition was exciting and worthwhile. We presented in this paper a number of different views and analysis based on the data gathered during the first phase. Many new solvers were submitted this year, a majority of them were zchaff-like solvers. There was no outstanding improvements compared to last year, but a reasonable set of strong solvers for “industrial” benchmarks is now available. Forklift, which inherits from berkmin62, leads those solvers. In the handmade category, lsat demonstrated an interesting behavior (a former version participated to SAT’02 but was incorrect). Satzoo showed better performances in terms of diversity of benchmarks and scalability. In the random category, the warded solvers are not a surprise: the solver kcufs in the complete category, the incomplete solver unitwalk in the satisfiable category.

The most interesting result of the competition is the behavior of the portfolio algorithm satzilla, which compared favourably with the solvers awarded in the random category. Whether this result can be replicated in the other categories is an interesting challenge.

Next year, a new competition will be organized in conjunction with SAT2004, in Vancouver. We showed in this paper is that the competition is now mature enough to fulfill three main goals: (1) motivate the field, (2) promote new solvers and hard benchmarks and (3) learn from the observation of solvers behaviors.

## Acknowledgements

We especially want to thanks our three judges, John Franco, Hans van Maaren and Toby Walsh for their involvement in every phases of the competition. Authors would also like to thanks the “Laboratoire de Recherche en Informatique”

(LRI, Orsay, France) and the “Dipartimento di Informatica Sistemica e Telematica” (DIST Genoa, Italy) for providing us with clusters of machines. At last, we thank all the authors of solvers and benchmarks for their participation and their effort to make their solver I/O compliant with our format. Without them there would be no competition at all!

## References

1. *Sixth International Conference on Theory and Applications of Satisfiability Testing*, S. Margherita Ligure - Portofino ( Italy), May 2003.
2. Nadel Alexander. Backtrack search algorithms for propositional satisfiability: Review and innovations. Master’s thesis, Hebrew University of Jerusalem, November 2002. please email alikn@hotmail.com to receive an electronic copy of this work.
3. Gilles Audemard, Daniel Le Berre, Olivier Roussel, Inês Lynce, and Jo ao Marques Silva. OpenSAT: An Open Source SAT Software Project. In *Sixth International Conference on Theory and Applications of Satisfiability Testing* [1].
4. A. Biere, A. Cimatti, E.M. Clarke, O. Strichman, and Y. Zhu. *Highly Dependable Software* , volume 58 of *Advances In Computers*, chapter Bounded model checking. Academic Press, 2003.
5. Armin Biere, Alessandro Cimatti, Edmund M. Clarke, M. Fujita, and Y. Zhu. Symbolic model checking using SAT procedures instead of bdds. In *Proceedings of Design Automation Conference (DAC’99)*, 1999.
6. Koen Claessen and Niklas Srensson. Finite model generation for first order logic using propositional satisfiability. Submitted for publication.
7. Gilles Dequen and Olivier Dubois. Renormalization as a function of clause lengths for solving random k-sat formulae. In *Proceedings of Fifth International Symposium on Theory and Applications of Satisfiability Testing*, pages 130–132, 2002.
8. Olivier Dubois and Gilles Dequen. A backbone-search heuristic for efficient solving of hard 3-sat formulae. In *Proceedings of the Seventeenth International Joint Conference on Artificial Intelligence (IJCAI’01)*, Seattle, Washington, USA, August 4th-10th 2001.
9. Niklas Een and Niklas Srensson. Temporal induction by incremental sat solving. Submitted for publication.
10. E. Goldberg and Y. Novikov. BerkMin: A fast and robust SAT-solver. In *Design, Automation, and Test in Europe (DATE ’02)*, pages 142–149, March 2002.
11. E. A. Hirsch and A. Kojevnikov. UnitWalk: A new SAT solver that uses local search guided by unit clause elimination. PDMI preprint 9/2001, Steklov Institute of Mathematics at St.Petersburg, 2001. A journal version is submitted to this issue.
12. Henry A. Kautz and Bart Selman. Pushing the envelope : Planning, propositional logic, and stochastic search. In *Proceedings of the Twelfth National Conference on Artificial Intelligence (AAAI’96)*, pages 1194–1201, 1996.
13. O. Kullmann. Heuristics for SAT algorithms: Searching for some foundations, September 1998. 23 pages, updated September 1999 w.r.t. running times, url = <http://cs-svr1.swan.ac.uk/csoliver/heur2letter.ps.gz>.
14. Daniel Le Berre, Laurent Simon, and Armando Tachella. Challenges in the QBF arena: the SAT’03 evaluation of QBF solvers. In *Proceedings of SAT2003*, 2003. submitted.

15. K. Leyton-Brown, E. Nudelman, G. Andrew, J. McFadden, and Y. Shoham. A portfolio approach to algorithm selection. In *Proceedings of IJCAI'03*, 2003. A full version of this paper is available from <http://robotics.stanford.edu/~kevinlb/boosting.pdf>, and is under review.
16. K. Leyton-Brown, E. Nudelman, and Y. Shoham. Learning the empirical hardness of optimization problems: The case of combinatorial auctions. In *Proceedings of CP'02*, 2002.
17. X. Y. Li, M.F. Stallmann, and F. Brglez. QingTing: A Fast SAT Solver Using Local Search and Efficient Unit Propagation. In *Sixth International Conference on Theory and Applications of Satisfiability Testing* [1]. See also <http://www.cbl.ncsu.edu/publications/>, and <http://www.cbl.ncsu.edu/OpenExperiments/SAT/>.
18. Inês Lynce and Joao Marques Silva. On implementing more efficient data structures. In *Sixth International Conference on Theory and Applications of Satisfiability Testing* [1].
19. Joo P. Marques-Silva and Karem A. Sakallah. Boolean Satisfiability in Electronic Design Automation. In *Proceedings of the IEEE/ACM Design Automation Conference (DAC)*, pages 675–680, June 2000.
20. Fabio Massacci. Using walk-sat and rel-sat for cryptographic key search. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence (IJCAI'99)*, pages 290–295, Stockholm, Sweden, July 31-August 6 1999. Morgan Kaufmann.
21. M. W. Moskewicz, C. F. Madigan, Y. Zhao, L. Zhang, and S. Malik. Chaff: Engineering an efficient SAT solver. In *Proceedings of the 38th Design Automation Conference (DAC'01)*, pages 530–535, June 2001.
22. R. Ostrowski, E. Grégoire, B. Mazure, and L. Sais. Recovering and exploiting structural knowledge from cnf formulas. In *Proc. of the Eighth International Conference on Principles and Practice of Constraint Programming (CP'2002)*, LNCS, pages 185–199, Ithaca (N.Y.), September 2002. Springer.
23. S. D. Prestwich. Randomised backtracking for linear pseudo-boolean constraint problems. In *Proceedings of Fourth International Workshop on Integration of AI and OR techniques in Constraint Programming for Combinatorial Optimisation Problems*, 2002.
24. Bart Selman, Henry A. Kautz, and David A. McAllester. Ten challenges in propositional reasoning and search. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence (IJCAI'97)*, pages 50–54, 1997.
25. Laurent Simon and Philippe Chatalic. SATEX: a web-based framework for SAT experimentation. In Henry Kautz and Bart Selman, editors, *Electronic Notes in Discrete Mathematics*, volume 9. Elsevier Science Publishers, June 2001. <http://www.lri.fr/~simon/satex/satex.php3>.
26. Laurent Simon, Daniel Le Berre, and Edward E. Hirsch. The sat2002 competition report. *Annals of Mathematics and Artificial Intelligence*, 2003. Special issue for SAT2002, to appear.
27. Geoff Sutcliffe and Christian Suttner. Evaluating general purpose automated theorem proving systems. *Artificial Intelligence*, 131:39–54, 2001.
28. Hantao Zhang. SATO: an efficient propositional prover. In *Proceedings of the International Conference on Automated Deduction (CADE'97)*, volume 1249 of *LNAI*, pages 272–275, 1997.
29. L. Zhang, C. F. Madigan, M. W. Moskewicz, and S. Malik. Efficient conflict driven learning in a Boolean satisfiability solver. In *International Conference on Computer-Aided Design (ICCAD'01)*, pages 279–285, November 2001.