# PARALLEL SAT SOLVING

*Jan Eberhardt, Jakub Lichman*

Department of Computer Science
ETH Zurich
Zurich, Switzerland

## ABSTRACT

Satisfiability solving is an interesting problem with various applications. In this report we describe how we implemented and parallelized a simple SAT solver. We compare different parallelization variations and analyze their respective trade offs. We achieved linear or even super linear speedups until up to 28 CPU cores compared to our sequential implementation.

## 1. INTRODUCTION

**Motivation.** The boolean satisfiability problem (SAT) belongs to the most important problems in program analysis, verification and other disciplines of theoretical computer science. It is used in the background of many applications, especially ones in the field of automated planning and scheduling, model checking (formal verification) and theorem proving. The last decade brought many improvements to the SAT solving world in the form of advanced heuristics, pre-processing and processing techniques, and data structures that allow efficient implementation of search space pruning. The past 10 years were also rich in improvements in parallelism. Current trends in computer hardware design decreased performance per processing unit and pack more units on a single processor. However, algorithms such as DPLL and CDCL were invented before the wide use of parallelism and therefore designed for sequential execution. Since the SAT problem is NP-complete, we consider it a good candidate to be parallelized.

In our approach, we are trying to speed up SAT solving by running it on multiple cores with different techniques of search space partitioning.

The remainder of this report is structured as follows: Next subsection provides related work. In Section 2 we introduce core SAT concepts. Section 3 explains how we parallelized DPLL and section 4 contains experimental evaluation. In Sections 5, 6 and 7 we talk about experiments we tried that did not work out, about possible extensions to our work and finally conclusion.

**Related work.** Tomas Balyo et al. proposed the "Horde-Sat" solver, which can run on up to 1024 cores and is based on the CDCL algorithm. Their parallel approach is different from ours, because it is portfolio based but with signs of search space partitioning. [1]

Most of the previous SAT solvers designed for computer clusters or grids use explicit search space partitioning. Examples of such solvers are GridSAT [2], PMSAT [3] or ManySat [4].

Jurkowiak et al. proposed a work stealing implementation for dynamic load-balancing. [5] Our work stealing approach is similar to what they proposed.

Berger et al. parallelized DPLL in Haskell. [6] Our DPLL search space partitioning is similar to what they proposed, but our solver is not limited to one machine and we can run on multiple compute nodes.

## 2. BACKGROUND

**CNF.** A *boolean variable* is a variable that can be assigned either to *true* or *false*. A *literal* of a boolean variable $x$ is considered to be in positive $x$ or negative $\overline{x}$ form. A *clause* is then disjunction (OR) of literals. A *conjunctive normal form* (CNF) formula is a conjunction of such clauses. Every propositional logic formula can be transformed to a CNF format. For the sake of this report we assume that the formula is already in the CNF format. However a transformation would be simple and straight forward based on basic propositional logic equivalence rules.

**SAT Solver.** A SAT solver is a program that is able to decide whether given formula is satisfiable. More formally, a given formula $F$ is satisfiable *iff* there *exists* an assignment of literals $\theta$ that makes whole formula (*CNF*) true. Such a formula $F$ is then called *satisfiable*. If there does not exits such an assignment, the formula is *unsatisfiable*. Furthermore, a SAT solver should be able to provide a valid variable assignment in the *satisfiable* case.

**History.** The first SAT solving algorithm was develop in 1960 by Martin Davis and Hilary Putnam [7]. They used a resolution-based decision procedure. Since the Davis-Putnam (or DP) algorithm was able to handle just valid formulas, a more general approach for SAT solving was needed.
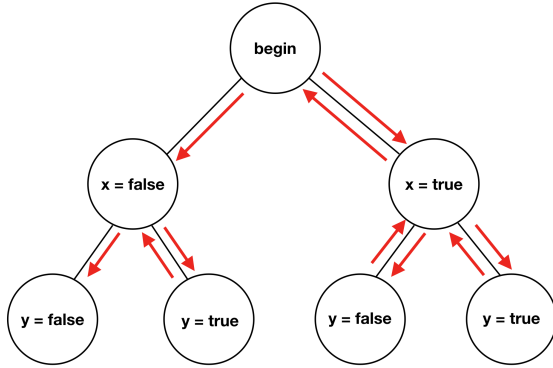
**Fig. 1**. The way how DPLL explores search space of CNF with two variables *x* and *y*.

//TODO: background check! I think DP can easily be extended to also handle unsat cases if I remember correctly...

In 1962 Davis, Putnam, Logemann and Loveland developed an extension of the DP algorithm, the DPLL algorithm. DPLL is a backtracking-based search algorithm that still forms the basis for many extensions and efficient implementations of SAT solvers.

**DPLL.** DPLL solves a SAT problem by modeling it as a decision tree, variables can either be assigned to *true* or *false*. In every step DPLL tries to find variables that are "automatically" assigned because of a previous decision, if there are none it will pick a variable $v_i$ and assume a value for it. If it later turns out that this decision was wrong it backtracks to that point and picks the negated assignment for variable $v_i$. Figure 1 shows how DPLL's decision tree could look like for a simple formula with 2 variables.

**CDCL.** The algorithm performs as well as DPLL a depth-first search of the space of partial truth assignments. In addition to DPLL, Conflict-Driven Clause Learning (*CDCL*) adopts a pruning technique called learning. While in DPLL it can happen that the same clause is solved multiple times, CDCL remembers them and in the next encounter avoids them. More formally, learning extracts and memorizes information from the previously searched space to prune the search for future branches [8]. Learning is done by adding clauses to the existing clause set. Clauses are analyzed and stored whenever the search reaches a conflict state. If it cannot be resolved by backtracking then the formula is unsatisfiable. If all variables are assigned and no conflict happened, then the formula is satisfiable.

## 3. PARALLELIZING DPLL

DPLL is a simple backtracking algorithm that explores the search space without any advanced techniques like learning and pruning. As a result of this, it does not require any trans-

fer of information between stages of solving and therefore it can be nicely parallelized. The parallelization is not trivial, but far easier than in the case of CDCL, where in addition to load balancing we need to deal with learned clause sharing.

**DPLL Branches.** As shown in Figure 1 the DPLL algorithm explores the search space with depth first search. It does that as follows: If there are no more variables that can be assigned trivially, one variable is assumed to be either *true* or *false*. That is exactly the point where we can do so called DPLL branching: continue on one branch and let some other node/processor solve the other branch. We can call this point a branching point and it represents one node (or decision) in the DPLL decision tree.

In our work we considered two approaches to load balancing. The first one is a master worker communication model, where one node, called *master*, is in charge of storing partial models and passing them to *workers* that will solve them. The second approach is a work stealing communication model, where each node (or *worker*) runs on its own and all *workers* are equal. If some *worker* runs out of work and the formula is not solved yet, it picks a random other *worker* and tries to steal a partial model from it, such that it can continue doing some work.

**Master Worker Model.** First we decided to design and implement the master worker model. We considered this model to be the easiest to implement and therefore we picked it first. The biggest challenge in parallelizing DPLL is load balancing. The work cannot be split equally between workers at the beginning, because it would require knowledge of the depth of each branch. This information cannot be inferred from a formula without solving it. We want to avoid situations where work is not divided equally because we want to fully exploit the computing power of all nodes.

The *master* is one of the *n* available nodes and the other *n-1* nodes are *workers*. At the beginning all workers are marked as free. The *master* starts solving by passing an empty model to the first free worker. Each worker in every branching point takes one branch and sends the other branch as a partial model to the master. When the master receives a partial model, it checks whether there exists some free worker. If not it stores the model into its local stack. Otherwise it sends it to a worker, which finished solving a branch with *unsat* and therefore is free and needs more work. This procedure is repeated until some worker finds a *sat* model and sends it to the master. The master then outputs the model and stops all workers. If no worker found a *sat* model, the master has an empty stack and all workers are free, then we know that the formula is *unsat*.

While the master worker model is easy to understand and implement, it comes with several drawbacks. The first of them is the lack of one node. In the master worker model, one node (the master) is only used as a manager and therefore this node is not actually doing any computation. The

next issue is related to scaling. The master can become a bottleneck with a certain number of nodes, because there is too much communication it needs to process and therefore workers need to wait longer for getting new partial models. The last problem is also related to the amount of communication in this model. Partial models and meta data are transferred every time a worker branches or wants to get a new partial model. This communication could be reduced. There is also communication at the beginning and at the end but this communication is necessary to distribute initial work and stop all workers.

**Work Stealing Model.** Because of the drawbacks mentioned above we have decided to reduce and decentralize the communication by implementing a work stealing model. The model treats all nodes *equally* and therefore there is *no master*, just workers. All workers contain their own stacks with partial models to be processed. This reduces the communication overhead because all produced models are stored locally instead of sending them to and storing them at the master. However, the existence of this local stack breaks load balancing because some workers can work on bigger subtrees of the search space than the others.

Inequalities in load balancing are resolved by a mechanism called work stealing. If some worker finished its work and no other worker has solved the formula yet, than it tries to "steal" a partial model from one of the other workers. The process of stealing can be defined as follows: A worker tries to steal from a randomly selected other worker until it finds a model or until someone else has found one. There are two more phases in the communication model. The first phase is the initial work distribution and the third one is the stopping criteria. In the first phase worker 0 will take the role of the master and distribute starting models to all other workers. Then all workers will try to solve their subtrees and in the case of running out of work they will try to steal work from some other worker (phase 2). When some worker finds a *sat* model, it will send it to worker 0, who stops all workers and outputs the final model (phase 3). A master node is necessary at the end because there can be multiple workers that find a *sat* model at the same time. In our implementation, worker 0 prints the first obtained *sat* model and ignores the others.

Our current work stealing implementation has one drawback: It can only handle satisfiable formulas. The stopping criteria in the unsatisfiable case was very trivial in the master worker model because there exists a moment when master's stack was empty and all workers were free. In the work stealing model, no such a moment exits. We do not have central node that controls the others and therefore knows in which stage each worker is and therefore we restrict ourselves just to satisfiable formulas with this communication model.

**Implementation.** We implemented both communica-

tion models in C++ with the help of MPI [9]. Our DPLL solver is implemented in C++ as well and therefore we can directly interact with MPI from within the solver. We tried to keep the implementation of the DPLL algorithm and communication within nodes as separate as possible for better readability. The implementation of parallel DPLL is straight forward and tries to follow good object oriented programming principles.

We used *vector*s from the C++ standard library to represent sets of clauses and sets of literals. We used *unsigned int*s as variable names, and *boolean* types for sign and value of the literals and variables. To implement the backtracking algorithm we used recursion, but we allocate all necessary data structures on the heap. We did not experience any stack overflow issues.

**Correctness.** We tested our sequential DPLL solver on hundreds of formulas that we either randomly generated or took over from a SATLIB formula collection. [10] We ran each formula through z3 [11] and compared the result of z3 with our result. The following cases have to be considered for each test case: If both solvers return unsat, we pass the test case. If one solver returns sat and the other one unsat, we fail the test case. If both solvers return sat, we still need to check if the model our solver has returned is correct. To do that we can conjoin the model to the original formula and again run it through z3.

We ran both parallel implementations through the same set of formulas and checked correctness for each of them. Assuming that our sequential implementation is correct, it is straight forward that our parallel implementation is correct as well: we essentially solved the same set of sub problems but on different nodes and in a potentially different order. The difference between the parallel and stealing version is just the load balancing and not in the searched space. Both parallel versions solve the same partial models but they are distributed over the working nodes in a different way.

## 4. EXPERIMENTAL RESULTS

In this section we are going to present our experimental results for both versions of our parallel DPLL algorithm as well as their detailed description.

**Experimental Setup.** We ran both communication models on the Euler super compute cluster [12] on up to 48 cores, which was the maximum accessible number of cores to us, and requested 1 Gigabyte of memory per core. Euler contains 5 different types of nodes but each of them contains at least two 12-core Intel Xeon processors (2.5-3.7 GHz) and between 64 and 512 GB of DDR4 memory clocked at 1866, 2133 or 2666 MHz.

**Benchmark Formulas.** During our correctness testing and debugging phase of the core algorithm, we realized that random formulas generated by our own random formula
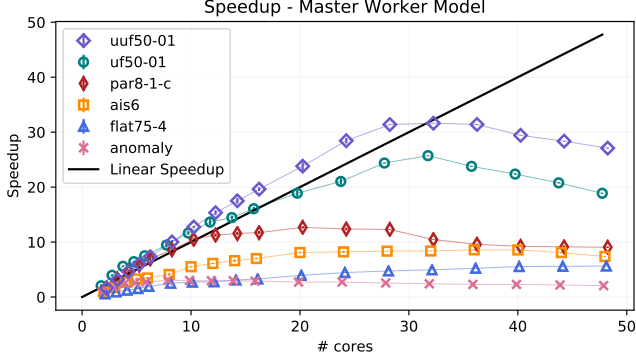
**Fig. 2**. Average speedup of master worker parallel DPLL implementation compared to sequential DPLL. The 95% confidence intervals are shown as error bars but too small to be visible in some cases.



**Fig. 3**. Average overall waiting time of workers per cnf in the parallel version. All waiting times per worker are summed up. The 95% confidence interval are shown as error bars.

generator are not suitable for performance testing. Even if we picked formulas that contained the same amount of variables, clauses and literals per clause, we observed large variations in run time between them. This is caused by the completely random structure of formulas generated by our generator. Therefore we used a set of 14 formulas taken from a SATLIB formula collection as a benchmark set. The set contains different real world problems modeled as a boolean formulas from various domains such as planning problems, all-interval series encodings, flat graph coloring, inductive inference, etc. Our benchmark set also contains a few random formulas that were taken from the SATLIB formula collection as well. Table 2 in the Appendix section contains a detailed listing of the used formulas together with their description. For the coming subsections we reduced the set of benchmark formulas to a representable subset of size 6 for readability reasons. We included the corresponding figures for the other 8 formulas in the Appendix section.

**Master Worker.** We compared the run times of the master worker communication pattern with the sequential version of DPLL. The achieved speedup is shown in Figure 2.

| Formula | seq. runtime [ms] | # vars/clauses |
|---|---|---|
| uf50-01 (sat) | 7'890.9 ± 44.31 | 50/218 |
| par8-1-c (sat) | 2'185.0 ± 32.88 | 64/254 |
| ais6 (sat) | 2'964.0 ± 41.54 | 61/581 |
| flat75-4 (sat) | 26'284.6 ± 172.65 | 225/840 |
| anomaly (sat) | 279.7 ± 6.74 | 48/261 |
| uuf50-01 (unsat) | 13'404.2 ± 102.94 | 50/218 |

**Table 1**. Overview of sequential runtime of benchmark subset formulas. The second column shows the average runtime with the 95 percentiles.

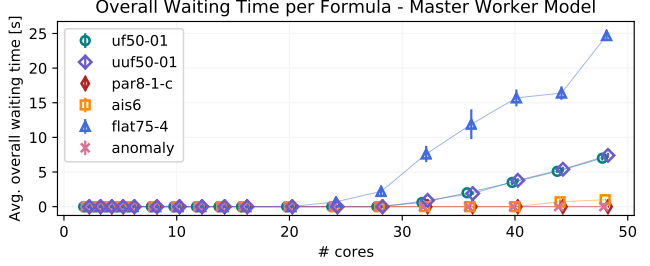It can be inferred from the figure that the achieved speedup

factor heavily depends on the formula. But it is not just the "size" of the problem or time required to solve a formula with the sequential algorithm that influences the speedup factor. As shown in Table 1 the formulas where we reached the highest speedup are not necessarily the ones with the longest time to solve sequentially or the largest number of variables (highest upper bound on the depth of the decision tree). Note that the best speedup in this subset is achieved for the randomly generated formula (uf50-01). For the random formulas it is easy to explain why the speedup goes down with more than 32 cores. It is caused by moments in the solving process when the master is no longer able to serve all the workers at the same time. Therefore some workers have to wait longer to obtain the next partial model. This problem starts to occur when the number of cores is larger than 30. The overall average waiting time of all workers, shown in Figure 3, shows this problem. The waiting time is considered to be the time that some worker spends waiting for a new partial model from the master. The overall waiting time is the sum of all waiting time periods of all workers, averaged over different runs of the same formula.

The cases where the speedup is not so significant are real world problems. The structure of their decision trees is less random. The small speedup is caused by the way the decision tree is traversed. While the sequential implementation deterministically does a depth-first search and therefore gives us some ordering guarantees, the parallel version sends all subbranches to the master and gives no ordering guarantees. It might happen that a worker is "close" to a solution but then receives a partial model in a completely different part of the decision tree that might be unsatisfiable all together. The size of the subtree where we will not find a solution depends on the structure of the formula. This is also the explanation why we achieved super linear speedup in some other cases: There we were lucky enough to reach the final solution with less iterations than in the sequential case.

**Work Stealing.** Similar to the previous comparison, we

compared our parallel work stealing model to the sequential DPLL version. The speedup that we achieved is shown in Figure 4. With the work stealing model we achieved large speedups for a small number of cores but essentially did not gain from adding more cores after some limit. The reason for this is again the way how we traverse the decision tree. In theory we should do something that is similar to a global breath first search and local (per node) depth first search in this model. However because of phase 0 mentioned in Section 3, the "global breath first" search part is not really a breath first search. Each worker starts in a node in the decision tree that is one level shifted inwards from the right-most subbranch of the tree. That means that if the formula is satisfiable, we are guaranteed that at least one worker is working on the correct subtree already after the first decision and will never solve models of the other subtree.

As a result, we outperform master worker with the work stealing version with only a few cores. But do not gain that much from adding more cores since the search space is split in a non-uniform way.

**Communication Overhead.** One of the reasons why we introduced a work stealing mechanism, was that it should drastically reduce the amount of communication (bytes transferred) between cores. Our experimental results prove that and they are shown in Figure 6. From the upper subplot in Figure 6 it can be inferred that the amount of communication is quite high, but does not necessarily increase when increasing the number of cores. This constant behavior is caused by the fact that every worker sends every partial model it finds and that it does not directly solve back to the master, regardless of the number of cores that are used. If more cores are used however, it can happen that parts of the decision tree that were not traversed with a smaller number of cores are suddenly traversed. In the lower subplot of Figure 6 the amount of overall communication is reduced, compared to the master worker model. This reduction is caused by the reduced amount of partial model transfers, as mentioned in Section 3. However as we discussed in the previous subsection, the amount of communication is not the only limiting factor. Master becoming the bottleneck when adding more cores is the main factor. The measured waiting time per worker did also not increase as much for the work stealing version as it did increase for the master worker model when adding more cores.

## 5. DISCUSSION

In this section we are going to discuss some approaches that we tried during development but that did not work as expected or we had to stop investigating possible solutions because of time constraints.

**DPLL-CDCL Hybrid Parallel Solver.** We have implemented a fully working CDCL solver. Its correctness was
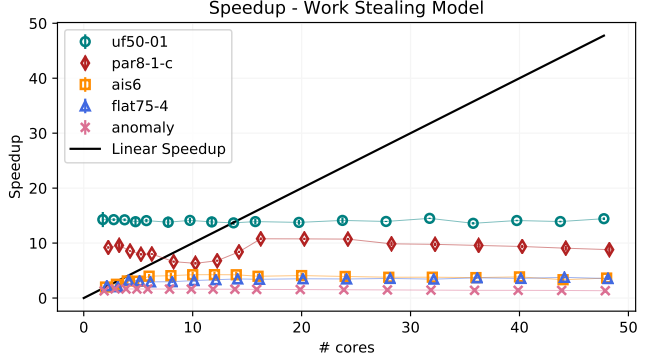


**Fig. 4**. Speedup of work stealing parallel DPLL implementation compared to sequential DPLL.
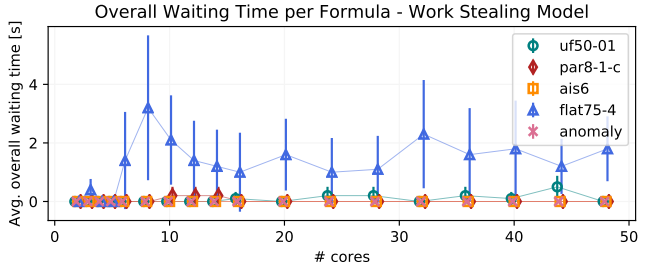


**Fig. 5**. Average overall waiting time of workers per cnf in the work stealing model. All waiting times per worker are summed up.
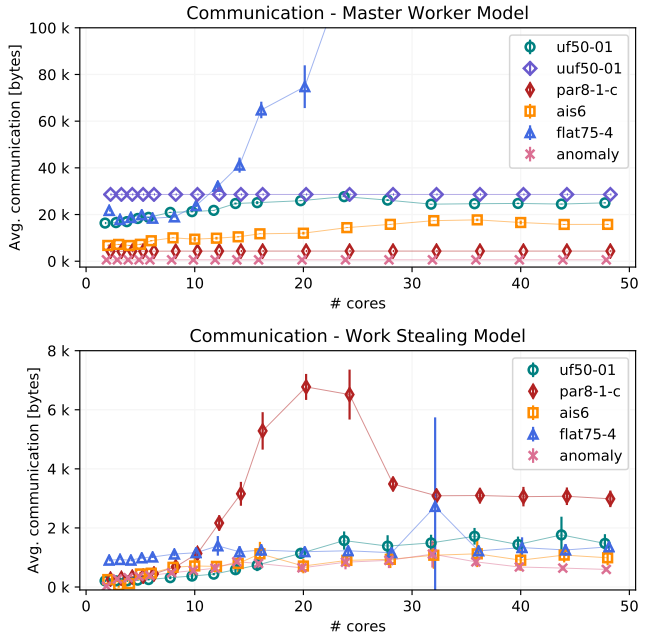


**Fig. 6**. Overall amount of transferred bytes for both the master worker model and the stealing model. The 95% confidence intervals are shown as error bars but too small to be visible for some data points.

proven with the same testing infrastructure that we use to test our DPLL implementation. The sequential version of the CDCL solver was faster than the sequential DPLL implementation for all of our benchmark formulas. Quite naturally we tried to plug that solver in locally at the workers to boost the performance. The way we did it is the following:

- First we run DPLL in parallel (it does not matter which of the two communication patterns we pick) and we only branch a certain number of times per worker.

- After that "branching limit" is reached we switch to CDCL and solve the whole subtree of the original decision tree locally.

- When CDCL terminates, we either found a solution or get a new model (either from the master or another worker) and solve it again with CDCL.

This hybrid approach showed worse performance than we expected. The reason for this is the following: When running DPLL globally, we essentially create multiple smaller problems. Smaller problems for the DPLL algorithm, which means they are easier and faster to solve than the original problem for a DPLL solver. For the CDCL algorithm on the other hand those sub-problems are not necessary easier. Some of those sub problems might actually be a lot harder than the original problem, because with the made assumption by a DPLL decision step we could add a new conflict that was not there in the original problem.

We ran everything we discussed in Section 4 with the hybrid parallel solver for different branching limits. We do not include any further information in this report because of space limitations. Overall the performance of the DPLL-CDCL hybrid solver was worse than both sequential CDCL and parallel DPLL with both of the two parallelization models.

## 6. FUTURE WORK

In this section we are going to discuss possible extensions and improvements of our work.

**Improvements.** As presented in Section 4 during the evaluation and measurements phase of our project we detected several drawbacks in our approach that could be improved.

First and the most important one is the decision tree traversal. One would need to ensure a more uniform distribution of the search space among the workers for both models. For instance we should avoid situations where one or no worker is working on the left part of the decision tree and all the other workers on the right part.

A second improvement would be to enhance the stealing effectiveness. Stealing should no longer be completely random. One could try to remember the stealing attempts

and steal from workers that have the most partial models in their local stacks.

**Extensions.** Our current work stealing model implementation is only able to handle *sat* formulas. For handling *unsat* formulas one would need to detect if all workers are trying to steal work.

Another interesting extension would be to fully parallelize the CDCL algorithm. It would require sharing of learned clauses, which could be implemented by extending both of our communication models.

## 7. CONCLUSION

SAT solving is after six decades of intensive research still a hot topic in theoretical computer science. In our approach we tried to leverage more processor cores to get more computing power for solving SAT problems. We parallelized the DPLL algorithm with two different communication models. A master worker model with a central management node (the master) and a work stealing model where we try to reduce the communication overhead to a minimum. Even though we identified deficiencies in both of our models we were still able to outperform our sequential DPLL implementation with both models. For some formulas we were even able to outperform our sequential CDCL implementation. We also investigated in a hybrid approach with global DPLL and local CDCL, which unfortunately did not bring the desired performance increase.

## 8. REFERENCES

[1] Tomas Balyo, Peter Sanders, and Carsten Sinz, "Hordesat: A massively parallel portfolio sat solver," *Theory and Applications of Satisfiability Testing*, vol. 18, pp. 156–172, 2015.

[2] W. Chrabakh and R. Wolski, "chaff-based distributed sat solver for the grid.," *ACM/IEEE conference on Supercomputing.*, vol. 37, 2003.

[3] L. Gil, P. Flores, and Silveira L.M., "Pmsat: a parallel version of minisat.," *Journal on Satisfiability, Boolean Modeling and Computation*, vol. 6, pp. 73–101, 2008.

[4] Y. Hamadi, S. Jabbour, and L. Sais, "Manysat: a parallel sat solver.," *Journal on Sat- isfiability, Boolean Modeling and Computation*, vol. 6, pp. 245262, 2008.

[5] Bernard Jurkowiak, Chu Min, and LiGil Utard, "A parallelization scheme based on work stealing for a class of sat solvers," *Journal of Automated Reasoning*, vol. 34, pp. 73101, 2005.

[6] Till Berger and David Sabel, "Parallelizing dpll in haskell," 02 2013.

[7] Martin Davis and Hilary Putnam, "A computing procedure for quantification theory," *J. ACM*, vol. 7, no. 3, pp. 201–215, July 1960.

[8] Lintao Zhang, Conor F. Madigan, Matthew H. Moskewicz, and Sharad Malik, "Efficient conflict driven learning in a boolean satisfiability solver," in *Proceedings of the 2001 IEEE/ACM International Conference on Computer-aided Design*, Piscataway, NJ, USA, 2001, ICCAD '01, pp. 279–285, IEEE Press.

[9] "MPI Forum," online: http://mpi-forum.org/, Accessed: 15.12.2017.

[10] "SATLIB - Benchmark Problems, University of British Columbia," online: http://www.cs.ubc.ca/~hoos/SATLIB/benchm.html, Accessed: 15.12.2017.

[11] Leonardo De Moura and Nikolaj Bjørner, "Z3: An efficient smt solver," in *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, Berlin, Heidelberg, 2008, TACAS'08/ETAPS'08, pp. 337–340, Springer-Verlag.

[12] "Euler High-Performance-Cluster, ETH Zurich," online: https://scicomp.ethz.ch/wiki/Euler, Accessed: 15.12.2017.

## 9. APPENDIX

This section contains additional information that is not strictly part of the report.

| File names | Description |
|---|---|
| anomaly, medium | SAT-encoded blocks world planning problems |
| ais6, ais8 | SAT-encoded All-Interval Series problems |
| flat50-1, flat75-4 flat 75-8 | SAT-encoded "Flat" Graph coloring problems (J. Coberson's flat graph generator is used |
| ii8a1 | Inductive inference, stem from a formulation of boolean function synthesis problems |
| par8-1-c, par8-4- | Instances for learning the paritiy function |
| uf50-01, uf75-01 | Uniform random 3-sat, satisfiable |
| uuf50-01, uuf50-02 | Uniform random 3-sat, unsatisfiable |

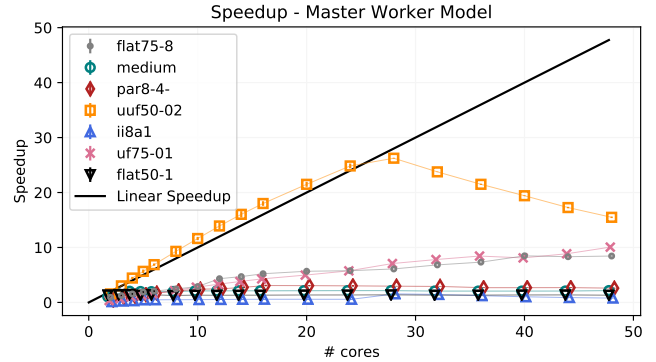**Table 2**. Description of benchmark problems.



**Fig. 7**. Average speedup of master worker parallel DPLL implementation compared to sequential DPLL. The 95% confidence intervalls are shown as error bars but too small to be visible in some cases. Second subset.
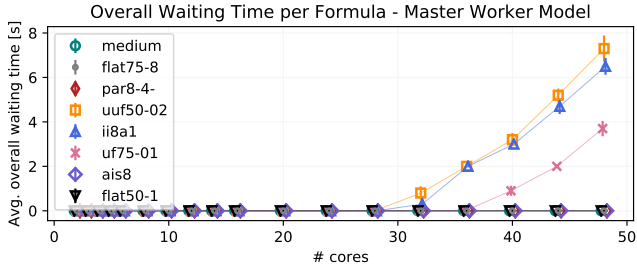
**Fig. 8**. Average overall waiting time of workers per cnf in parallel version. All waiting times per worker are summed up. The 95% confidence interval are shown as error bars. Second subset.
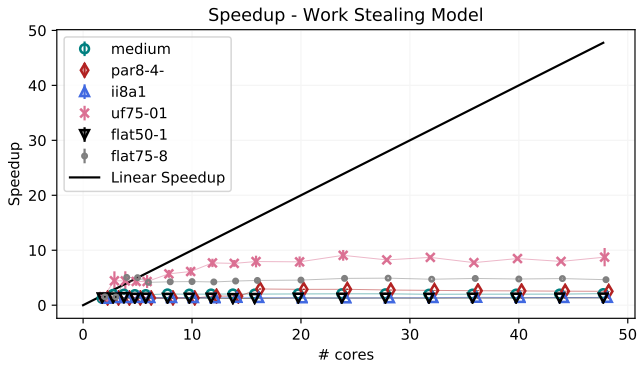


**Fig. 9**. Speedup of work stealing parallel DPLL implementation compared to sequential DPLL. Second subset.

| Formula | Best Solver Configuration | Runtime [ms] |
|---------|---------------------------|--------------|
| anomaly | CDCL sequential | $16.2 \pm 4.18$ |
| medium | CDCL sequential | $414.8 \pm 32.94$ |
| ais6 | DPLL master worker, 14 cores | $343.3 \pm 3.58$ |
| ais8 | DPLL master worker, 48 cores | $3015.3 \pm 58.85$ |
| flat50-1 | CDCL sequential | $36.4 \pm 9.07$ |
| flat75-4 | CDCL sequential | $702.8 \pm 38.18$ |
| flat75-8 | CDCL sequential | $1054.2 \pm 50.72$ |
| ii8a1 | DPLL master worker, 6 cores | $72.2 \pm 0.79$ |
| par8-1-c | CDCL sequential | $15.2 \pm 2.28$ |
| par8-4- | DPLL master salve, 2 cores | $132.2 \pm 2.22$ |
| uf50-01 | CDCL sequential | $31.5 \pm 9.21$ |
| uf75-01 | DPLL master worker, 32 nodes | $297.4 \pm 7.99$ |
| uuf50-01 | DPLL master worker 24 nodes | $187.5 \pm 5.59$ |
| uuf50-02 | DPLL master worker 24 nodes | $148.9 \pm 2.85$ |

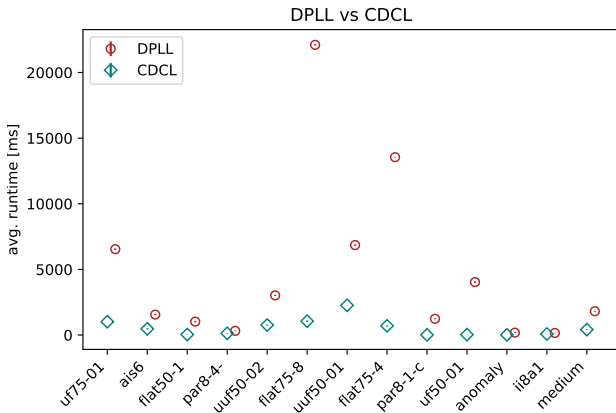**Table 3**. Benchmark formulas with best performing configuration.



**Fig. 10**. DPLL vs CDCL sequential runtime. Average runtime with 95% confidence intervals.