

# Verifying Solvers: How Much do You Want to Prove?

---

Mathias Fleury

January 17, 2024



# How do We Make SAT Solvers Correct?

## Proofs

- requires to check the proof for each file
- not all techniques can be represented by current proof formats

Program Verification This talk!

- works for every input, so no overhead
- does not crash even if run the program for a year

# How do We Make SAT Solvers Correct?

## Proofs

- requires to check the proof for each file
- not all techniques can be represented by current proof formats

## Program Verification This talk!

- works for every input, so no overhead
- does not crash even if run the program for a year

# Restrictions

- Parsing is always trusted `CakeML` has some modelisation of file systems, but don't try `grep aaa file >> file` at home
- Printing the answer is trusted too
- The resulting SAT solvers live outside of their system
  - you cannot use them in the system you do the proofs

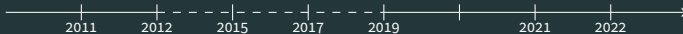
# A Personal History of Solver Verification

Lescuyer, CDCL, PhD thesis, Coq  
Maric and Janicic, LMCS, Isabelle, SML code  
CDCL+2WL,

Berger, Lawrence, Forsberg,  
and Seisenberger, DPLL, Minlog,  
Haskell code

Fleury, Blanchette, Weidenbach,  
CDCL+2WL, Isabelle, SML code

Fleury (and Lammich), CDCL+2WL,  
Isabelle, EDA Challenge, LLVM IR  
code



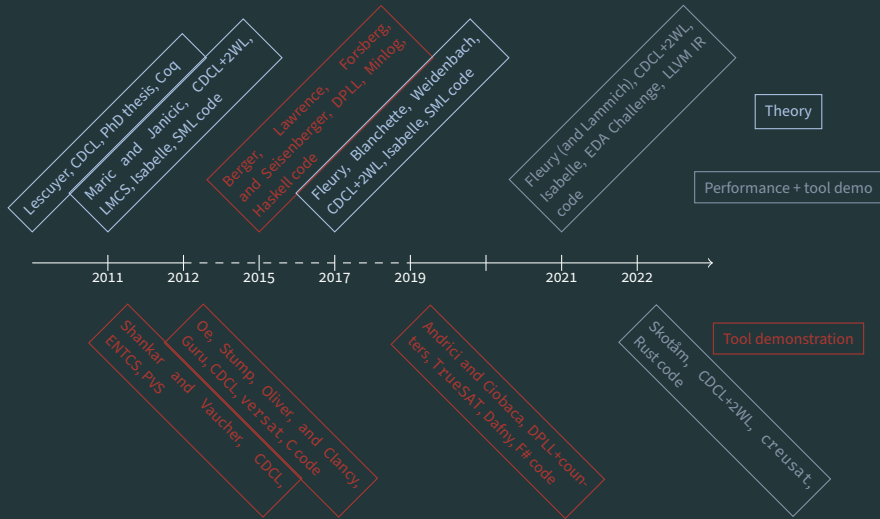
Shankar and Vaucher, ENTCS, PVS  
Oe, Stump, Oliver, and Clancy,  
Guru, CDCL, versat, C code  
CDCL,

Andric and Ciobaca, DPLL+counters, TrueSAT, Dafny, F# code

Skoták, CDCL+2WL, creusat,  
Rust code

incomplete especially because the bottom-up approach is a good master thesis

# A Personal History of Solver Verification



## What Far Can You Go With One Master Thesis?

**Fleury:** functional code, DPLL, no restarts, propagation by going over all clauses, decision by going over all clauses. *but it terminates and is complete*  
Solves no problem from the SAT Comp

**Skotåm:** imperative code, CDCL, restarts, watch literals, decision heuristic.  
Solves > 150.

**Top:** Some theory expressed in your tool

**Bottom** Some (hopefully fast) code

All full verifications go top-down.

seL4 kernel is mixed:

Specification  $\rightarrow$  Haskell  $\leftarrow$  C

Most partial verifications go bottom-up.

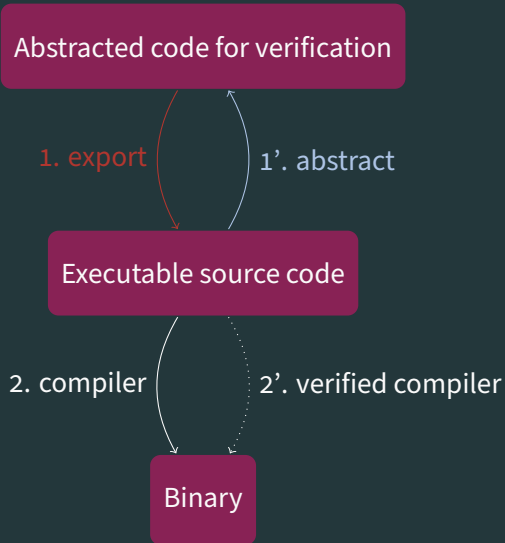
Most natural for each tool!



# Bottom-Up Or the Art of Proving very little

---

# Organisation



Translation from Rust to why3 (unverified) [Denis, Jourdan, Marché, ICFEM'21]

1' translation from Guru to C [Stump et al, PLPV'09]

2': only used in a SAT checker

**Implicit Checker** The checker = the verification

Every approach I am aware of: checker = resolution checker

**Theorem (Correctness)**

*Deriving  $\perp$  implies that the problem is UNSAT.*

**Implicit Checker** The checker = the verification

Every approach I am aware of: checker = resolution checker

Theorem (Correctness)

*Deriving  $\perp$  implies that the problem is UNSAT.*

**Implicit Checker** The checker = the verification

Every approach I am aware of: checker = resolution checker

## Theorem (Correctness)

*Deriving  $\perp$  implies that the problem is UNSAT.*

## Some Invariants of a SAT Solver

**Deriving the empty clause:** input problem unsat

**Conflicts on current level:** runtime assertion

**Termination:** Unknown

**No conflict+all assigned:** checking of the model

**No crash:** depends on the approach

# What Do You Have To Prove?

**Well-behaved:** no read past end of array

**Clauses:** not modified except by resolution

- But: non trivial for minimization where the resolution is implicit

## Making the Solver more Complex: Adding Restart?

Assume you already have a working CDCL.

Adding restarts means:

1. call backtrack to level 0

That is all

except for heuristics, performance debugging, ...



# Challenges

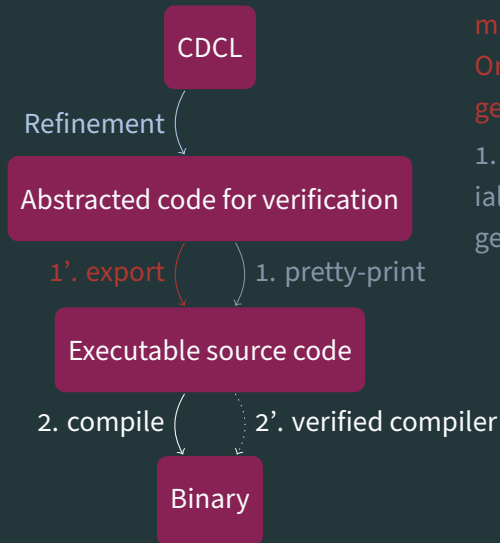
## What is hard?

- Usually relies on automatic provers, which must be able to handle the specification Skotåm: swapping literals
- No termination ITP don't like non-termination
- Closer to programs written by hand easier to try different strategies

## Top-Down Approach: Proving Too much

---

# Organisation



1'. Hupel: use semantics from 2'.  
Or Lammich: LLVM generation

1. trusted as trivial translation (SML generation)

## Key Idea

### Abstract Correctness (Pragmatic) CDCL is fully correct

#### Theorem (Total Correctness<sup>1</sup>)

*Deriving  $\perp$  iff the problem is UNSAT. No conflict + total assignment = SAT. Termination.*

#### Theorem (Total Correctness IsaSAT-LLVM)

*If the answer is not unknown, it is either SAT with a model or UNSAT.*

IsaSAT-SML had full correction SML semantics does not forbid arrays  $\geq 2^{64}$ ,  
no compiler support

---

<sup>1</sup>At some point, memory representation can cause also aborts.

## Key Idea

**Abstract Correctness** (Pragmatic) CDCL is fully correct

### Theorem (Total Correctness<sup>1</sup>)

*Deriving  $\perp$  iff the problem is UNSAT. No conflict + total assignment = SAT. Termination.*

### Theorem (Total Correctness IsaSAT-LLVM)

*If the answer is not unknown, it is either SAT with a model or UNSAT.*

IsaSAT-SML had full correction SML semantics does not forbid arrays  $\geq 2^{64}$ ,  
no compiler support

---

<sup>1</sup>At some point, memory representation can cause also aborts.

## Key Idea

**Abstract Correctness** (Pragmatic) CDCL is fully correct

### Theorem (Total Correctness<sup>1</sup>)

*Deriving  $\perp$  iff the problem is UNSAT. No conflict + total assignment = SAT. Termination.*

### Theorem (Total Correctness IsaSAT-LLVM)

*If the answer is not unknown, it is either SAT with a model or UNSAT.*

IsaSAT-SML had full correction – SML semantics does not forbid arrays  $\geq 2^{64}$ ,  
no compiler support

---

<sup>1</sup>At some point, memory representation can cause also aborts.

# Key Idea

**Abstract Correctness** (Pragmatic) CDCL is fully correct

## Theorem (Total Correctness<sup>1</sup>)

*Deriving  $\perp$  iff the problem is UNSAT. No conflict + total assignment = SAT. Termination.*

## Theorem (Total Correctness IsaSAT-LLVM)

*If the answer is not unknown, it is either SAT with a model or UNSAT.*

IsaSAT-SML had full correction SML semantics does not forbid arrays  $\geq 2^{64}$ ,  
no compiler support

---

<sup>1</sup>At some point, memory representation can cause also aborts.

# Refinement in IsaSAT





# Some Invariants of a SAT Solver

**Deriving the empty clause:** unsat (OR: derive conflict at level 0)

**Conflicts on current level:** completeness of propagations

**Termination:** Yes

IsaSAT can answer unknown

if too many clauses  $\sum_{c \in \text{clauses}} 5 + |c| \approx |\text{clause\_memory}| \geq 2^{63}$

**No crash:** yes (up to the assumptions on memory) allocation does not fail

## Making the Solver more Complex: Adding Restart?

Assume you already have a working CDCL.

Adding restarts means:

1. change your CDCL (to include a counter to increase restart interval)
2. change the refinement to be based on the extended CDCL
3. add restarts with the counter. Make sure that it does not overflow.

That is all

except for heuristics, performance debugging, ...

# Challenges

## What is hard?

- you have to prove everything [lots of code](#)
- limited by the speed of your tools [bring Isabelle to its knees](#)
- hard to find people [Isabelle and code synthesis can be seen as two different systems](#)

# Refinements

In retrospect over the entire project:

- Many components that are not independent everything is parametrized by the set of variables...  
Watch list can be indexed by every literal in the set of clauses
- Mistakes have been made: too much coupling ... that is not duplicated  
Better: watch lists are defined over a set of literals that is the same as the set of clauses
- But: refactoring takes time.

# Refinements

In retrospect over the entire project:

- Many components that are not independent everything is  
parametrized by the set of variables...

Watch list can be indexed by every literal in the set of clauses

- Mistakes have been made: too much coupling ... that is not duplicated

Better: watch lists are defined over a set of literals that is the same as the set of clauses

- But: refactoring takes time.

# Refinements

In retrospect over the entire project:

- Many components that are not independent everything is parametrized by the set of variables...

Watch list can be indexed by every literal in the set of clauses

- Mistakes have been made: too much coupling ... that is not duplicated

Better: watch lists are defined over a set of literals that is the same as the set of clauses also moving up proof that index valid

- But: refactoring takes time.

# Refinements

In retrospect over the entire project:

- Many components that are not independent everything is parametrized by the set of variables...

Watch list can be indexed by every literal in the set of clauses

- Mistakes have been made: too much coupling ... that is not duplicated

Better: watch lists are defined over a set of literals that is the same as the set of clauses also moving up proof that index valid

- But: refactoring takes time.

# Refinements

In retrospect over the entire project:

- Many components that are not independent everything is parametrized by the set of variables...  
Watch list can be indexed by every literal in the set of clauses
- Mistakes have been made: too much coupling ... that is not duplicated  
Better: watch lists are defined over a set of literals that is the same as the set of clauses also moving up proof that index valid
- But: refactoring takes time.



# Refinements

In retrospect over the entire project:

- Many components that are not independent everything is parametrized by the set of variables...  
Watch list can be indexed by every literal in the set of clauses
- Mistakes have been made: too much coupling ... that is not duplicated  
Better: watch lists are defined over a set of literals that is the same as the set of clauses also moving up proof that index valid
- But: refactoring takes time.

# Refinements

In retrospect over the entire project:

- Testing new features hard Some I implemented and proved things that did not work and I removed.
- Testing improvement for code generation structure was forced, not a choice. Pointers

## What Can You Not Express?

- aliasing

```
struct ISASAT {  
    TRAIL trail;  
    CLAUSES clauses; ....  
};
```

```
struct ISASAT solver;  
isasat->trail = assign(lit, solver->trail);
```

- pointers are complicated    IsaSAT: I tried to use a pointer to a state and never managed to make it less than 10 times slower

## What Can You Not Express?

- aliasing

```
struct ISASAT {  
    TRAIL trail;  
    CLAUSES clauses; ....  
};
```

```
struct ISASAT solver;  
isasat->trail = assign(lit, solver->trail);
```

- pointers are complicated    IsaSAT: I tried to use a pointer to a state and never managed to make it less than 10 times slower

# The Code

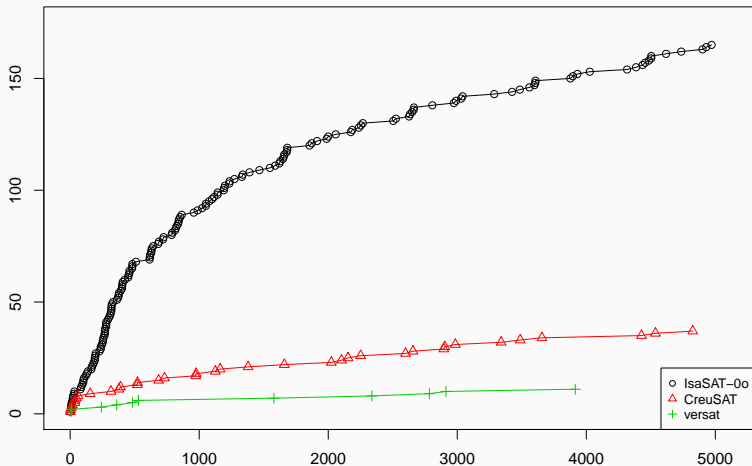
---

## How Do They Perform?

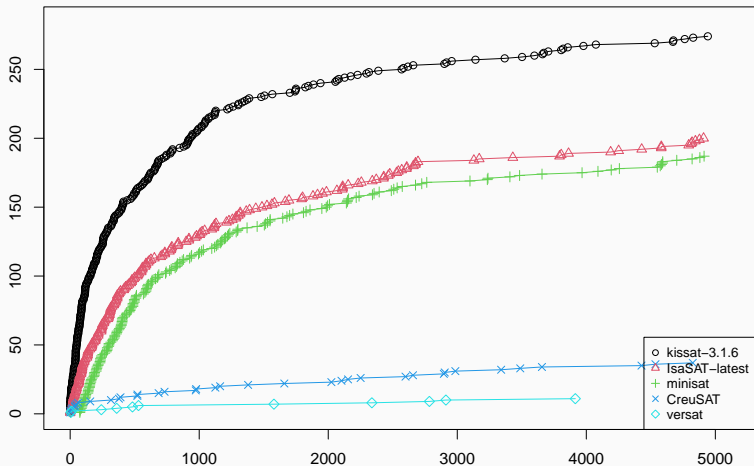
Solver	SAT	UNSAT
IsaSAT	175	130
Creusat	145	79
versat	60	62

**Table 1:** Results on the SC2015 according to Skotåm (24 GB, 1800 s)

# How Do They Perform?

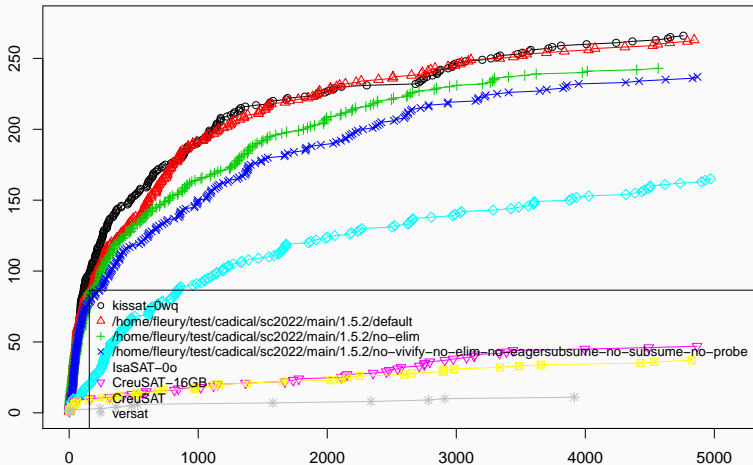


# How Do They Perform?





# How Do They Perform?



## How Good Is The Code (I)? Guru

```
void * gpropagate_h(int gnv_24, int gdl_4, void * gas_37, void * gws_17) {
    start_gpropagate_h: {
        /* match with exactly one case: gassign_state */
        void * gpa_13;
        void * gwhy_6;
        void * gdls_6;
        void * ghist_6;
        int ghist_cur_4;
        int ghist_end_4;
        void * gcarraway_tmp_119;
        gpa_13 = ginit_unique_unique(guarray, gas_37, ((gAssignState_gassign_state *)gas_37)->gpa_2)
        ;
        gwhy_6 = ginit_unique_unique(gwarray, gas_37, ((gAssignState_gassign_state *)gas_37)->gwhy_2)
        ;
        gdls_6 = ginit_unique_unique(guarray, gas_37, ((gAssignState_gassign_state *)gas_37)->gdls_2)
        ;
        ghist_6 = ginit_unique_unique(guarray, gas_37, ((gAssignState_gassign_state *)gas_37)->ghist_2)
        ;
        [...]
        switch ((int)gcarraway_tmp_120) {

        case op_gff: {

            fprintf(stderr, "abort at /Users/kain/Projects/versat/old_versions/0.6/src/unitprop.g, line 76");
        }
    }
}
```

## How Good is The Code (II)? IsaSAT

(If times permits)

HTML version of the Isabelle files: [https://people.mpi-inf.mpg.de/~mfleury/IsaFoL/current/Weidenbach\\_Book/IsaSAT/IsaSAT\\_Inner\\_Propagation\\_Defs.html#IsaSAT\\_Inner\\_Propagation\\_Defs.unit\\_propagation\\_update\\_statistics|const](https://people.mpi-inf.mpg.de/~mfleury/IsaFoL/current/Weidenbach_Book/IsaSAT/IsaSAT_Inner_Propagation_Defs.html#IsaSAT_Inner_Propagation_Defs.unit_propagation_update_statistics|const)

Correcntess theorem: [https://people.mpi-inf.mpg.de/~mfleury/IsaFoL/current/Weidenbach\\_Book/IsaSAT/IsaSAT\\_All\\_LLVM.html](https://people.mpi-inf.mpg.de/~mfleury/IsaFoL/current/Weidenbach_Book/IsaSAT/IsaSAT_All_LLVM.html)

## How Good Is The Code (II)? IsaSAT

```
define ISASAT_STATE @unit_propagation_outer_loop_wl_D(ISASAT_STATE %x) #0 {  
  
  start:  
    %x1 = call i8 @IsaSAT_Profile_PROPAGATE ()  
    call void @IsaSAT_Profile_LLVM_start_profile (i8 %x1)  
    br label %while_start  
  
  while_start:  
    %s = phi ISASAT_STATE [ %x3, %while_body ], [ %x, %start ]  
    %x2 = call i1 @literals_to_update_wl_empty_fast_code (ISASAT_STATE %s)  
    br i1 %x2, label %while_body, label %while_end  
  
  while_body:  
    %xb = call { ISASAT_STATE, i32 } @select_and_remove_from_literals_to_update_wl(ISASAT_STATE %s)  
    %a1 = extractvalue { ISASAT_STATE, i32 } %xb, 0  
    %a2 = extractvalue { ISASAT_STATE, i32 } %xb, 1  
    %x3 = call ISASAT_STATE @unit_propagation_inner_loop_wl_D (i32 %a2, ISASAT_STATE %a1)  
    br label %while_start  
  
  while_end:  
    %xc = call i8 @IsaSAT_Profile_PROPAGATE ()  
    call void @IsaSAT_Profile_LLVM_stop_profile (i8 %xc)  
    ret ISASAT_STATE %s  
}
```

(only edit: ISASAT\_STATE is unfolded in the code and remove prefix from function names)

## How Good Is The Code (II)? CreuSAT

```
#[cfg_attr(feature = "trust_unit", trusted)]  
#[ensures(f.equisat(^f))]  
pub fn unit_propagate(f: &mut Formula, trail: &mut Trail, watches: &mut Watches) -> Result<()  
    let mut i = trail.curr_i;  
    let old_trail: Ghost<&mut Trail> = ghost! { trail };  
    let old_f: Ghost<&mut Formula> = ghost! { f };  
    let old_w: Ghost<&mut Watches> = ghost! { watches };  
    #[invariant(trail_inv, trail.invariant(*f))]  
    while i < trail.trail.len() {  
        let lit = trail.trail[i].lit;  
        match propagate_literal(f, trail, watches, lit) {  
            Ok(_) => {}  
            Err(cref) => {  
                return Err(cref);  
            }  
        }  
        i += 1;  
    }  
    trail.curr_i = i;  
    Ok(())  
}
```

(only edit: remove some invariants and ensures)

# Conclusion

---

## Comparison: How different are there really?

- Removing assertions from bottom-up means being more top-down and requires more proofs where automation struggles
- Very hard to remove proofs from top-down
- Link top-down with concrete code? Currently has not been tried but I am trying to find a student

## Conclusion

- Only application of verified SAT solvers: finishing last at SAT Competition, getting Masters, or PhDs
- Unexpectedly, IsaSAT correlates the least with Kissat on SC2022 benchmarks observation by Max Heisinger
- But: do you have applications where proof checking is not possible?
- What is the right timeout for the SAT Competition? In 2023: 5000 s solving +  $9 \times 5000$  s checking. with 5x in total, probably still behind Laurent, but not last