

A Formally Verified Checker for First-Order Proofs

Soulkee Baek

3rd February, 2021

Abstract

The Verified TESC Verifier (VTV) is a formally verified checker for the new Theory-Extensible Sequent Calculus (TESC) proof format for first-order ATPs. VTV accepts a TPTP problem and a TESC proof as input, and uses the latter to verify the unsatisfiability of the former. VTV is written in Agda, and the soundness of its proof-checking kernel is verified in respect to a first-order semantics formalized in Agda. VTV shows robust performance in a comprehensive test using all eligible problems from the TPTP problem library, successfully verifying all but the largest 5 of 12296 proofs, with >97% of the proofs verified in less than 1 second.

1 Introduction

Modern automated reasoning tools are highly complex pieces of software, and it is generally no simple matter to establish their correctness. Bugs have been discovered in them in the past [18, 11], and more are presumably hidden in systems used today. One popular strategy for coping with the possibility of errors in automated reasoning is the *De Bruijn* criterion [2], which states that automated reasoning software should produce ‘proof objects’ which can be independently verified by simple checkers that users can easily understand and trust. In addition to reducing the trust base of theorem provers, the De Bruijn criterion comes with the additional benefit that the small trusted core is often simple enough to be formally verified itself. Such thoroughgoing verification is far from universal, but there has been notable progress toward

this goal in various subfields of automated reasoning, including interactive theorem provers, SAT solvers, and SMT solvers.

One area in which similar developments have been conspicuously absent is first-order automated theorem provers (ATPs), where the lack of a machine-checkable proof format [17] precluded the use of simple independent verifiers. The Theory-Extensible Sequent Calculus (TESC) is a new proof format for first-order ATPs designed to fill this gap. In particular, the format’s small set of fine-grained inference rules makes it relatively easy to implement and verify its proof checker.

This paper presents the Verified TESC Verifier (VTV), a proof checker for the TESC format ¹ written and verified in Agda [3]. The aim of the paper is twofold. Its primary purpose is to demonstrate the reliability of TESC proofs by showing precisely what is established by their successful verification using VTV. Its secondary aim is to serve as a guide for understanding VTV’s codebase and its design choices, which can be especially useful for readers who want to implement their own TESC verifiers.

The rest of the paper is organized as follows: Section 2 gives a brief survey of similar works and how VTV relates to them. Section 3 describes the syntax and inference rules of the TESC proof calculus. Section 4 presents the main TESC verifier kernel, and Section 5 gives a detailed specification of the verifier’s soundness property. Sections 3, 4, and 5 also include code excerpts and discuss how their respective contents are formalized in Agda. Section 6 shows the results of empirical tests measuring VTV’s performance. Section 7 gives a summary and touches on potential future work.

2 Related Works

SAT solving is arguably the most developed subfield of automated reasoning in terms of verified proof checkers. A non-exhaustive list of SAT proof formats with verified checkers include RAT [13], RUP and IORUP [14], LRAT [6], and GRIT [7]. In the related field of SMT solving, the SMTCoq project [1] also uses a proof checker implemented and verified in the Coq proof assistant.

Despite the limitations imposed by Gödel’s second incompleteness theorem [10], there has been interesting work toward verification of interactive theorem provers. All of HOL Light except the axiom of infinity has been

¹<https://github.com/skbaek/tesc>

proven consistent in HOL Light itself [11], which was further extended later to include definitional principles for new types and constants [15]. More recent projects go beyond the operational semantics of programming languages and verifies theorem provers closer to the hardware, such as the Milawa theorem prover which runs on a Lisp runtime verified against x86 machine code [8], and the Metamath Zero [4] theorem prover verified down to x84-64 instruction set architecture.

VTV is designed to serve a role similar to these verified checkers for first-order ATPs and the TESC format. There has been several different approaches [21, 5] to verifying the output of first-order ATPs, but the only example I am aware of which uses independent proof objects with a verified checker is the Ivy system [16]. The largest difference between Ivy and VTV is that Ivy’s proof objects only record resolution and paramodulation steps, so all input problems have to be normalized by a separate preprocessor written in ACL2. In contrast, the TESC format supports preprocessing steps like Skolemization and CNF normalization, so it allows ATPs to work directly on input problems with their optimized preprocessing strategies.

3 Proof Calculus

The syntax of the TESC calculus is as follows:

$$\begin{aligned}
f &::= \sigma \mid \#_k \\
t &::= x_k \mid f(\vec{t}) \\
\vec{t} &::= \cdot \mid \vec{t}, t \\
\phi &::= \top \mid \perp \mid f(\vec{t}) \mid \neg\phi \mid \phi \vee \chi \mid \phi \wedge \chi \mid \phi \rightarrow \chi \mid \phi \leftrightarrow \chi \mid \forall\phi \mid \exists\phi \\
\Gamma &::= \cdot \mid \Gamma, \phi
\end{aligned}$$

f ranges over *functors*, which usually called ‘non-logical symbols’ in other presentations of first-order logic. The TESC calculus makes no distinction between function and relation symbols, and relies on the context to determine whether a symbol applied to arguments is a term or an atomic formula. For brevity, we borrow the umbrella term ‘functor’ from the TPTP syntax and use it to refer to any non-logical symbol.

There are two types of functors: σ ranges over *plain functors*, which you can think of as the usual relation or function symbols. We assume that there is a suitable set of symbols Σ , and let $\sigma \in \Sigma$. Symbols of the form $\#_k$

are *indexed functors*, and the number k is called the *functor index* of $\#_k$. Indexed functors are used to reduce the cost of introducing fresh functors: if you keep track of the largest functor index k that occurs in the environment, you may safely use $\#_{k+1}$ as a fresh functor without costly searches over a large number of terms and formulas.²

We use t to range over terms, \vec{t} over lists of terms, ϕ over formulas, and Γ over sequents. Quantified formulas are written without variables thanks to the use of De Bruijn indices [9]; the number k in variable x_k is its De Bruijn index. As usual, parentheses may be inserted for scope disambiguation, and the empty list operator ‘.’ may be omitted when it appears as part of a complex expression. E.g., the sequent \cdot, ϕ, ψ and term $f(\cdot)$ may be abbreviated to ϕ, ψ and f .

Formalization of TESC syntax in Agda is mostly straightforward, but with one caveat. A first instinct definition of terms might be something like the following:

```
data Term : Set where
  var : Nat → Term
  fun : List Term → Term
```

But this definition works poorly in practice, since any function defined by structural recursion on terms runs into non-termination issues. For instance, substitution on terms must make recursive calls with arguments of functors (the members of the list argument to `fun`), but Agda can’t automatically ascertain the termination of these calls because those terms are behind a `List` and hence not immediate subterms. Although it is possible to manually prove termination, it is much simpler to sidestep this issue with a pseudo-mutual recursion. The basic idea is to use a special type `Term* : Bool → Set` that can double as both terms and lists of terms depending on its boolean argument. I.e., we use `Term* false` as the type of terms, and `Term* true` as the type of lists of terms.

```
data Term* : Bool → Set where
  var : Nat → Term* false
  fun : Functor → Term* true → Term* false
  nil : Term* true
  cons : Term* false → Term* true → Term* true
```

²Thanks to Marijn Heule for suggesting this idea.

Rule	Conditions	Example
$\frac{\Gamma, A(b, \Gamma[i])}{\Gamma} A$		$\frac{\phi \leftrightarrow \psi, \phi \rightarrow \psi}{\phi \leftrightarrow \psi} A$
$\frac{\Gamma, B(0, \Gamma[i]) \quad \Gamma, B(1, \Gamma[i])}{\Gamma} B$		$\frac{\phi \vee \psi, \phi \quad \phi \vee \psi, \psi}{\phi \vee \psi} B$
$\frac{\Gamma, C(t, \Gamma[i])}{\Gamma} C$	$\text{lfi}(t) \leq \text{size}(\Gamma)$	$\frac{\neg \exists f(x_0), \neg f(g)}{\neg \exists f(x_0)} C$
$\frac{\Gamma, D(\text{size}(\Gamma), \Gamma[i])}{\Gamma} D$		$\frac{\exists f(x_0), f(\#_1)}{\exists f(x_0)} D$
$\frac{\Gamma, N(\Gamma[i])}{\Gamma} N$		$\frac{\neg \neg \phi, \phi}{\neg \neg \phi} N$
$\frac{\Gamma, \neg \phi \quad \Gamma, \phi}{\Gamma} S$	$\text{lfi}(\phi) \leq \text{size}(\Gamma)$	$\frac{\neg f(\#_0) \quad f(\#_0)}{\cdot} S$
$\frac{\Gamma, \phi}{\Gamma} T$	$\text{lfi}(\phi) \leq \text{size}(\Gamma),$ $\text{adm}(\text{size}(\Gamma), \phi)$	$\frac{= (f, f)}{\cdot} T$
$\frac{}{\Gamma} X$	For some i and j , $\Gamma[i] = \neg \Gamma[j]$	$\frac{}{\neg \phi, \phi} X$

Table 1: TESC inference rules.

which lets us define terms and lists of terms as

Term = **Term*** **false**
Terms = **Term*** **true**

The inference rules of the TESC calculus are shown in Table 1. The A, B, C, D , and N rules are the *analytic* rules which analyze existing formulas on the sequent and adds resulting subformulas to the new sequent (we are reading the proof tree in the bottom-up direction). The formula analysis functions used by rules are given in Table 2. Notice that the analytic rules are very similar to Smullyan's *uniform notation* for analytic tableaux, which is where they get their names from. Note that:

- $\Gamma[i]$ denotes the (0-based) i th formula of sequent Γ , where $\Gamma[i] = \top$ if the index i is out-of-bounds.

A	$A(0, \neg(\phi \vee \psi)) = \neg\phi$	$A(1, \neg(\phi \vee \psi)) = \neg\psi$
	$A(0, \phi \wedge \psi) = \phi$	$A(1, \phi \wedge \psi) = \psi$
	$A(0, \neg(\phi \rightarrow \psi)) = \phi$	$A(1, \neg(\phi \rightarrow \psi)) = \neg\psi$
	$A(0, \phi \leftrightarrow \psi) = \phi \rightarrow \psi$	$A(1, \phi \leftrightarrow \psi) = \psi \rightarrow \phi$
B	$B(0, \phi \vee \psi) = \phi$	$B(1, \phi \vee \psi) = \psi$
	$B(0, \neg(\phi \wedge \psi)) = \neg\phi$	$B(1, \neg(\phi \wedge \psi)) = \neg\psi$
	$B(0, \phi \rightarrow \psi) = \neg\phi$	$B(1, \phi \rightarrow \psi) = \psi$
	$B(0, \neg(\phi \leftrightarrow \psi)) = \neg\phi \rightarrow \psi$	$B(1, \neg(\phi \leftrightarrow \psi)) = \neg\psi \rightarrow \phi$
C	$C(t, \forall\phi) = \phi[0 \mapsto t]$	$C(t, \neg\exists\phi) = \neg\phi[0 \mapsto t]$
D	$D(k, \exists\phi) = \phi[0 \mapsto \#_k(\cdot)]$	$D(k, \neg\forall\phi) = \neg\phi[0 \mapsto \#_k(\cdot)]$

Table 2: Formula analysis functions. For any arguments not explicitly defined above, the functions all return \top . E.g., $A(0, \phi \vee \psi) = \top$.

- $\text{lfi}(x)$ denotes the largest functor index (lfi) occurring in x . If x includes no functor indices, $\text{lfi}(x) = -1$.
- $\text{size}(\Gamma)$ is the number of formulas in Γ .
- All inference rules are designed to preserve the invariant $\text{lfi}(\Gamma) < \text{size}(\Gamma)$ for every sequent Γ . We say that a sequent Γ is *good* if it satisfies this invariant.
- S is the usual cut rule, and X is the axiom or init rule.
- $\text{adm}(k, \phi)$ asserts that ϕ is an *admissible* formula in respect to the target theory and sequent size k . More precisely, if T is the target theory, Γ is a sequent satisfiable modulo T , and $\text{size}(\Gamma) = k$, then $\text{adm}(k, \phi)$ implies that adding ϕ to Γ preserves satisfiability modulo T . The sequent size argument k is required because some admissible formulas use this number as the functor index of newly introduced indexed functors. The T rule may be used to add any admissible formula.

The last part implies that the definition of the admissability predicate, and by extension the definition of well-formed TESC proofs, changes according to the implicit target theory. This is the ‘theory-extensible’ part of TESC. The current version of VTV verifies basic TESC proofs that target the theory

of equality, so it allows T rules to introduce equality axioms, fresh relation symbol definitions, and choice axioms. But VTV can be easily modified in a modular way to target other theories as well.

TESC proofs are formalized in Agda as follows:

```
data Proof : Set where
  rule-a : Nat → Bool → Proof → Proof
  rule-b : Nat → Proof → Proof → Proof
  rule-c : Nat → Term → Proof → Proof
  rule-d : Nat → Proof → Proof
  rule-n : Nat → Proof → Proof
  rule-s : Formula → Proof → Proof → Proof
  rule-t : Formula → Proof → Proof
  rule-x : Nat → Nat → Proof
```

There are parts of TESC proofs omitted in the definition of `Proof`, e.g. sequents. This is a design choice made in favor of efficient space usage. Since proofs are uniquely determined by their root sequents + complete information of the inference rules used, TESC proof files save space by omitting any components that can be constructed on the fly during verification, which includes all intermediate sequents and formulas introduced by analytic rules. Terms of the type `Proof` are constructed by parsing input TESC files, so it only includes information stored in TESC files, which are the arguments to the constructors of `Proof`.

4 The Verifier

Since `Proof` only includes basic information regarding inference rule applications, the verifier function for `Proof` must construct intermediate sequents as it recurses down a proof, and also check that inference rule arguments (e.g., the term t of a C rule application) satisfy their side conditions. This is exactly what `verify` does:

```
verify : Sequent → Proof → Bool
verify Γ (rule-a i b p) = verify (add Γ (analyze-a b (nth i Γ))) p
verify Γ (rule-b i p q) =
  (verify (add Γ (analyze-b false (nth i Γ))) p) ∧
  (verify (add Γ (analyze-b true (nth i Γ))) q)
```

```

verify  $\Gamma$  (rule-c  $i$   $t$   $p$ ) =
  term*-lfi<? (suc (size  $\Gamma$ ))  $t$   $\wedge$ 
  verify (add  $\Gamma$  (analyze-c  $t$  (nth  $i$   $\Gamma$ )))  $p$ 
verify  $\Gamma$  (rule-d  $i$   $p$ ) = verify (add  $\Gamma$  (analyze-d (size  $\Gamma$ ) (nth  $i$   $\Gamma$ )))  $p$ 
verify  $\Gamma$  (rule-n  $i$   $p$ ) = verify (add  $\Gamma$  (analyze-n (nth  $i$   $\Gamma$ )))  $p$ 
verify  $\Gamma$  (rule-s  $\phi$   $p$   $q$ ) =
  formula-lfi<? (suc (size  $\Gamma$ ))  $\phi$   $\wedge$ 
  verify (add  $\Gamma$  (not  $\phi$ ))  $p$   $\wedge$  verify (add  $\Gamma$   $\phi$ )  $q$ 
verify  $\Gamma$  (rule-t  $\phi$   $p$ ) =
  adm? (size  $\Gamma$ )  $\phi$   $\wedge$  verify (add  $\Gamma$   $\phi$ )  $p$ 
verify  $\Gamma$  (rule-x  $i$   $j$ ) = formula=? (nth  $i$   $\Gamma$ ) (not (nth  $j$   $\Gamma$ ))

```

Analytic rules introduce new formulas obtained by formula analysis using `analyze` functions, and side conditions are checked using appropriate `?` functions. The argument type `Sequent`, however, offers some interesting design choices. What kind of data structures should be used to encode sequents? The first version of VTV used lists, but lists immediately become a bottleneck with practically-sized problems due to their poor random access speeds. The default TESC verifier included in the TPTP-TSTP-TESC Processor (T3P, the main tool for generating and verifying TESC files) uses arrays, but arrays are hard to come by and even more difficult to reason about in dependently typed languages like Agda. Self-balancing trees like AVL or red-black trees come somewhere between lists and arrays in terms of convenience and performance, but it can still be tedious to prove basic facts about them if those proofs are not available in your language of choice, as is the case for Agda's standard library.

For VTV, we cut corners by taking advantage of the fact that (1) formulas are never deleted from sequents, and (2) new formulas are always added to the end of sequents. This allows us to use a simple tree structure:

```

data Tree (A : Set) : Set where
  empty : Tree A
  fork : Nat  $\rightarrow$  A  $\rightarrow$  Tree A  $\rightarrow$  Tree A  $\rightarrow$  Tree A

```

For any tree `fork k a t s` , the number k is the size of the left subtree t . This property is not guaranteed to hold by construction, but it is easy to ensure that it always holds in practice. With this definition, balanced addition of elements to trees becomes trivial:


```

add : {A : Set} → Tree A → A → Tree A
add empty a = fork 1 a empty empty
add ts@(fork k b t s) a =
  if (size s <b size t)
  then (fork (k + 1) b t (add s a))
  else (fork (k + 1) a ts empty)

```

Then the type `Sequent` can be defined as `Tree Formula`.

5 Soundness

In order to verify the soundness of `verify`, we first need to formalize a first-order semantics that it can be sound in respect to. Most of the formalization is routine, but it also includes some oddities particular to VTV.

One awkward issue that recurs in formalization of first-order semantics is the handling of arities. Given that each functor has a unique arity, what do you do with ill-formed terms and atomic formulas with the wrong number of arguments? You must either tweak the syntax definition to preclude such possibilities, or deal with ill-formed terms and formulas as edge cases, both of which can lead to unpleasant bloat.

For VTV, we avoid this issue by assuming that every functor has infinite arities. Or rather, for each functor f with arity k , there are an infinite number of other functors that share the name f and has arities $0, 1, \dots, k-1, k+1, k+2, \dots$ ad infinitum. With this assumption, the denotation of functors can be simply defined as

```

Rels : Set
Rels = List D → Bool

Funs : Set
Funs = List D → D

```

A `Rels` (resp. `Funs`) can be thought of as a collection of an infinite number of relations (resp. functions), one for each arity. An interpretation is a pair of a relation assignment and a function assignment, which assign `Rels` and `Funs` to functors.

```

RA : Set
RA = Functor → Rels

```

```

FA : Set
FA = Functor → Funs

```

variable assignments simply assign denotations to `Nat`, since variables are identified by their Bruijn indices.

```

VA : Set
VA = Nat → D

```

Term valuation requires a bit of tweaking due to the unusual definition of `Term*`:

```

ElemList : Set → Bool → Set
ElemList A false = A
ElemList A true = List A

term*-val : FA → VA → {b : Bool} → Term* b → ElemList D b
term*-val _ V (var k) = V k
term*-val F V (fun f ts) = F f (term*-val F V ts)
term*-val F V nil = []
term*-val F V (cons t ts) = (term*-val F V t) :: (term*-val F V ts)

```

Formula valuation is mostly straightforward, but some care needs to be taken in the handling of variable assignments. The truth value of universally quantified formulas are defined as:

$$R, F, V \models (\text{qtf false } \phi) = \forall x \rightarrow (R, F, (V / 0 \mapsto x) \models \phi)$$

`qtf` : `Bool` → `Formula` → `Formula` is a constructor of `Formula` for quantified formulas, and `qtf false` encodes the universal quantifier. For any V , k , and d , $(V / k \mapsto d)$ is an updated variable assignment obtained from V by assigning d to the variable x_k , and pushing the assignments of all variables larger than x_k by one. E.g., $(V / 0 \mapsto x) 0 = x$ and $(V / 0 \mapsto x) (\text{suc } m) = V m$.

Now we can define (un)satisfiability of sequents in terms of formula valuations:

```

satisfies : RA → FA → VA → Sequent → Set
satisfies R F V Γ = ∀ φ → φ ∈ Γ → R, F, V ⊨ φ

```

```

sat : Sequent → Set
sat Γ = ∃ λ R → ∃ λ F → ∃ λ V → (respects-eq R × satisfies R F V Γ)

unsat : Sequent → Set
unsat Γ = ¬ (sat Γ)

```

The `respects-eq` R clause asserts that the relation assignment R respects equality. This condition is necessary because we are targeting first-order logic with equality, and we are only interested in interpretations that satisfy all equality axioms.

VTV's formalization of first-order semantics is atypical in that (1) every functor doubles as both relation and function symbols with infinite arities, and (2) the definition of satisfiability involves variable assignments, thereby applying to open as well as closed formulas. But this is completely harmless for our purposes: whenever a traditional interpretation (with unique arities for each functor and no variable assignment) M satisfies a set of sentences Γ , M can be easily extended to an interpretation in the above sense that still satisfies Γ , since the truths of sentences in Γ are unaffected by functors or variables that do not occur in them. Therefore, if a set of sentences is unsatisfiable in the sense we've defined above, it is also unsatisfiable in the usual sense.

Now we finally come to the soundness statement for `verify`:

```

verify-sound : ∀ (S : Sequent) (p : Proof) →
  good S → T (verify S p) → unsat S

```

`T : Bool → Set` maps `true` to `T` and `false` to `⊥`. The condition `good S` is necessary, because the soundness of TESC proofs is dependent on the invariant that all sequents are good. But we can do better than merely assuming that the input sequent is good, because the parser which converts the input character list into the initial (i.e., root) sequent is designed to fail if the parsed sequent is not good. `parse-verify` is the outer function which accepts two character lists as argument, parses them into a `Sequent` and a `Proof`, and calls `verify` on them. The soundness statement for `parse-verify` is as follows:

```

parse-verify-sound : ∀ (seq-chars prf-chars : Chars) →
  T (parse-verify seq-chars prf-chars) →
  ∃ λ (S : Sequent) → returns parse-sequent seq-chars S × unsat S

```

In other words, if `parse-verify` succeeds on two input character lists, then the sequent parser successfully parses the first character list into some unsatisfiable sequent S . `parse-verify-sound` is an improvement over `verify-sound`, but it also shows the limitation of the current setup. It only asserts that there is *some* unsatisfiable sequent parsed from the input characters, but provides no guarantees that the parsed sequent is actually equivalent to the original TPTP problem. This means that the formal verification of VTV is limited to the soundness of its proof-checking kernel, and the correctness its TPTP parsing phase has to be taken on faith.

6 Test Results

The performance of VTV was tested by running it on all eligible problems in the TPTP [22] problem library. A TPTP problem is eligible if it satisfies all of the following conditions (parenthesized numbers indicate the total number of problems that satisfy all of the preceding conditions).

- It is in the CNF or FOF language (23291).
- Its status is ‘theorem’ or ‘unsatisfiable’ (13636).
- It conforms to the official TPTP syntax. More precisely, it does not have any occurrences of the character ‘%’ in the `sq_char` syntactic class, as required by the TPTP syntax. This is important because T3P assumes that the input TPTP problem is syntactically correct and uses ‘%’ as an endmarker (13389).
- All of its formulas have unique names. T3P requires this condition in order to unambiguously identify formulas by their names during proof compilation (13119).
- It can be solved by Vampire [19] or E [20] in one minute using default settings (Vampire = 7885, E = 4853 ³).

³It should be noted that the default setting is actually not optimal for E; when used with the `--auto` or `--auto-schedule` options, E’s success rates are comparable to that of Vampire. But the TSTP solutions produced by E under these higher-performance modes are much more difficult to compile down to a TESC proof, so they could not be used for these tests.

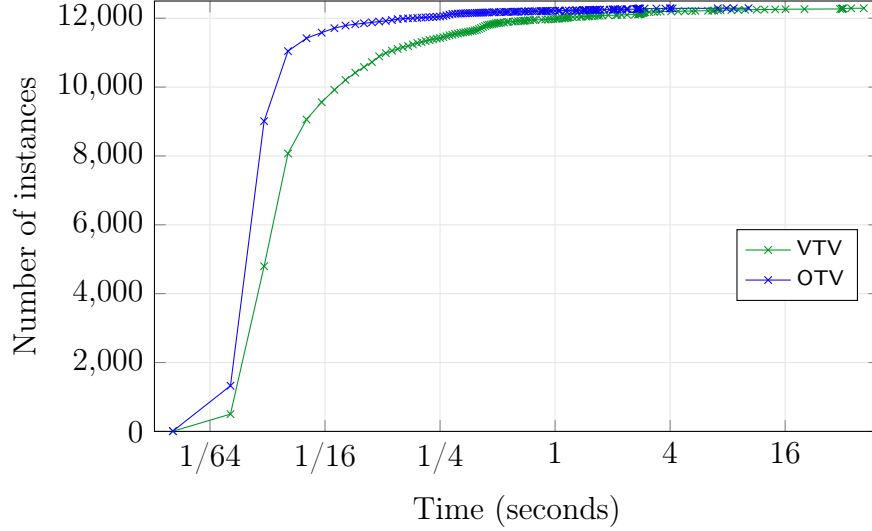


Figure 1: Verification times of VTV and OTV. The datapoints show the number of TESC proofs that each verifier could check within the given time limit. The plots look more “jumpy” toward the lower end due to the limited time measurement resolution (10 ms) of the unix `time` command.

- The TSTP solution produced by Vampire or E can be compiled to a TESC proof by T3P (Vampire = 7792, E = 4504).

The resulting $7792 + 4504 = 12296$ proofs were used for testing VTV. All tests were performed on Amazon EC2 `r5a.large` instances, running Ubuntu Server 20.04 LTS on 2.5 GHz AMD EPYC 7000 processors with 16 GB RAM.⁴

Out of the 12296 proofs, there were 5 proofs that VTV failed to verify due to exhausting the 16 GB available memory. A cactus plot of verification times for the remaining 12291 proofs are shown in Fig. 1. As a reference point, we also show the plot for the default TESC verifier included T3P running on the same proofs. The default TESC verifier is written in Rust, and is optimized for performance with no regard to verification. For convenience, we refer to it as the Optimized TESC Verifier (OTV).

VTV is slower than OTV as expected, but the difference is unlikely to be

⁴More information on the testing setup used, including the versions of software used, can be found at <https://github.com/skbaek/tesc/tree/itp>.

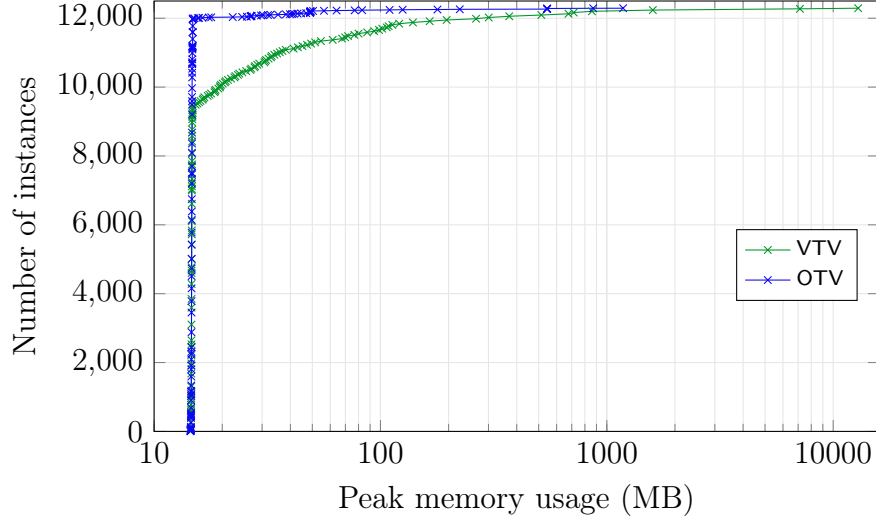


Figure 2: Peak memory usages of VTV and OTV. The datapoints show the number of TESC proofs that each verifier could check within the given peak memory usage.

noticed in actual use since the absolute times for most proofs are very short, and the total times are dominated by a few outliers. VTV verified $>97.4\%$ of proofs under 1 second, and $>99.3\%$ under 5 seconds. Also, the median time for VTV is 40 ms, a mere 10 ms behind the OTV’s 30 ms. But OTV’s mean time (54.54 ms) is still much shorter than that of VTV (218.93 ms), so users may prefer to use it for verifying one of the hard outliers or processing a large batch of proofs at once.

The main drawback of VTV is its high memory consumption. Fig. 2 shows the peak memory usages of the two verifiers. For a large majority of proofs, memory usages of both verifiers are stable and stay within the 14-20 MB range, but VTV’s memory usage spikes earlier and higher than OTV. Due the limit of the system used, memory usage could only be measured up to 16 GB, but the actual peak for VTV would be higher if we included the 5 failed verifications. A separate test running VTV on an EC2 instance with 64 GB ram (`r5a.2xlarge`) still failed for 3 of the 5 problematic proofs, so the memory requirement for verifying all 12296 proofs with VTV is at least >64 GB. In contrast, OTV could verify all 12296 proofs with less than 3.2 GB of memory.

7 Conclusion

VTV can serve as a fallback option whenever extra rigour is required in verification, thereby increasing our confidence in the correctness of TESC proofs. It also helps the design of other TESC verifiers by providing a reference implementation that is guaranteed to be correct. It should be particularly helpful for implementing other verified TESC verifiers in, say, Lean or Coq, since many of the issues we’ve discussed (termination checking, proving properties of trees, etc.) are common to these languages.

There are two main ways in which VTV could be further improved. Curb-ing its memory usage would be the most important prerequisite for making it the default verifier in T3P. This may require porting VTV to a verified programming language with finer low-level control over memory usage.

VTV could also benefit from a more reliable TPTP parser. A formally verified parser would be ideal, but the complexity of TPTP’s syntax makes it difficult to even *specify* the correctness of a parser, let alone prove it. A more realistic approach would be imitating the technique used by verified LRAT checkers [12], making VTV print the parsed problem and textually comparing its with the original problem.

Acknowledgements

This work has been partially supported by AFOSR grant FA9550-18-1-0120.

Thanks to Jeremy Avigad for a careful proofreading of drafts and advice for improvements.

References

- [1] Michaël Armand, Germain Faure, Benjamin Grégoire, Chantal Keller, Laurent Théry, and Benjamin Werner. A modular integration of sat/smt solvers to coq through proof witnesses. In *International Conference on Certified Programs and Proofs*, pages 135–150. Springer, 2011.
- [2] Henk Barendregt and Freek Wiedijk. The challenge of computer mathematics. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 363(1835):2351–2375, 2005.

- [3] Ana Bove, Peter Dybjer, and Ulf Norell. A brief overview of agda—a functional language with dependent types. In *International Conference on Theorem Proving in Higher Order Logics*, pages 73–78. Springer, 2009.
- [4] Mario Carneiro. Metamath zero: The cartesian theorem prover. *arXiv preprint arXiv:1910.10703*, 2019.
- [5] Zakaria Chihani, Tomer Libal, and Giselle Reis. The proof certifier checkers. In *International Conference on Automated Reasoning with Analytic Tableaux and Related Methods*, pages 201–210. Springer, 2015.
- [6] Luís Cruz-Filipe, Marijn JH Heule, Warren A Hunt, Matt Kaufmann, and Peter Schneider-Kamp. Efficient certified rat verification. In *International Conference on Automated Deduction*, pages 220–236. Springer, 2017.
- [7] Luís Cruz-Filipe, Joao Marques-Silva, and Peter Schneider-Kamp. Efficient certified resolution proof checking. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 118–135. Springer, 2017.
- [8] Jared Davis and Magnus O Myreen. The reflective milawa theorem prover is sound (down to the machine code that runs it). *Journal of Automated Reasoning*, 55(2):117–183, 2015.
- [9] Nicolaas Govert De Bruijn. Lambda calculus notation with nameless dummies, a tool for automatic formula manipulation, with application to the church-rosser theorem. In *Indagationes Mathematicae (Proceedings)*, volume 75, pages 381–392. North-Holland, 1972.
- [10] Kurt Gödel. Über formal unentscheidbare sätze der principia mathematica und verwandter systeme i. *Monatshefte für mathematik und physik*, 38(1):173–198, 1931.
- [11] John Harrison. Towards self-verification of hol light. In *International Joint Conference on Automated Reasoning*, pages 177–191. Springer, 2006.
- [12] Marijn Heule, Warren Hunt, Matt Kaufmann, and Nathan Wetzler. Ef-

- ficient, verified checking of propositional proofs. In *International Conference on Interactive Theorem Proving*, pages 269–284. Springer, 2017.
- [13] Marijn JH Heule, Warren A Hunt, and Nathan Wetzler. Verifying refutations with extended resolution. In *International Conference on Automated Deduction*, pages 345–359. Springer, 2013.
 - [14] Marijn JH Heule, Warren A Hunt Jr, and Nathan Wetzler. Bridging the gap between easy generation and efficient verification of unsatisfiability proofs. *Software Testing, Verification and Reliability*, 24(8):593–607, 2014.
 - [15] Ramana Kumar, Rob Arthan, Magnus O Myreen, and Scott Owens. Hol with definitions: Semantics, soundness, and a verified implementation. In *International Conference on Interactive Theorem Proving*, pages 308–324. Springer, 2014.
 - [16] William McCune and Olga Shumsky. Ivy: A preprocessor and proof checker for first-order logic. In *Computer-Aided reasoning*, pages 265–281. Springer, 2000.
 - [17] Giles Reger and Martin Suda. Checkable proofs for first-order theorem proving. In *ARCADE@ CADE*, pages 55–63, 2017.
 - [18] Giles Reger, Martin Suda, and Andrei Voronkov. Testing a saturation-based theorem prover: Experiences and challenges. In *International Conference on Tests and Proofs*, pages 152–161. Springer, 2017.
 - [19] Alexandre Riazanov and Andrei Voronkov. The design and implementation of vampire. *AI communications*, 15(2, 3):91–110, 2002.
 - [20] Stephan Schulz. E—a brainiac theorem prover. *Ai Communications*, 15(2, 3):111–126, 2002.
 - [21] Geoff Sutcliffe. Semantic derivation verification: Techniques and implementation. *International Journal on Artificial Intelligence Tools*, 15(06):1053–1070, 2006.
 - [22] Geoff Sutcliffe. The TPTP problem library and associated infrastructure. *Journal of Automated Reasoning*, 43(4):337, 2009.