

Formal Proofs of Tarjan's Strongly Connected Components Algorithm in Why3, Coq and Isabelle

Ran Chen

Institute of Software of the Chinese Academy of Sciences, Beijing, China
chenr@ios.ac.cn

Cyril Cohen

Université Côte d'Azur, Inria, Sophia Antipolis, France
cyril.cohen@inria.fr

Jean-Jacques Lévy

Irif & Inria, Paris, France
jean-jacques.levy@inria.fr

Stephan Merz

Université de Lorraine, CNRS, Inria, LORIA, Nancy, France
stephan.merz@inria.fr

Laurent Théry

Université Côte d'Azur, Inria, Sophia Antipolis, France
laurent.thery@inria.fr

Abstract

Comparing provers on a formalization of the same problem is always a valuable exercise. In this paper, we present the formal proof of correctness of a non-trivial algorithm from graph theory that was carried out in three proof assistants: Why3, Coq, and Isabelle.

2012 ACM Subject Classification Software and its engineering → Formal software verification

Keywords and phrases Mathematical logic, Formal proof, Graph algorithm, Program verification

Digital Object Identifier 10.4230/LIPIcs.ITP.2019.13

Supplement Material <http://www-sop.inria.fr/marelle/Tarjan/contributions.html>

1 Introduction

Graph algorithms are notoriously obscure in the sense that it is hard to grasp why exactly they work. Therefore, proofs of correctness are more than welcome in this domain. In this paper we consider Tarjan's algorithm [31] for computing the strongly connected components in a directed graph and present formal proofs of its correctness in three different systems: Why3, Coq and Isabelle/HOL. The algorithm is treated at an abstract level with a functional programming style manipulating finite sets, stacks and mappings, but it respects the linear time behaviour of the original presentation.

To our knowledge this is the first time that the formal correctness proof of a non-trivial program is carried out in three very different proof assistants: Why3 is based on a first-order logic with inductive predicates and automatic provers, Coq on an expressive theory of higher-order logic and dependent types, and Isabelle/HOL combines simply typed higher-order logic with automatic provers. Crucially for our comparison, the algorithm is defined at the same level of abstraction in all three systems, and the proof relies on the same arguments in the three formal systems. We deliberately decided not to base our representation of the algorithm on some specific infrastructure for program verification such as Lammich's monadic or imperative refinement frameworks [17, 19] for Isabelle/HOL or the existing encoding of Back and Morgan's refinement calculus in Coq [1], as doing so would make comparisons between



© Ran Chen, Cyril Cohen, Jean-Jacques Lévy, Stephan Merz, and Laurent Théry;
licensed under Creative Commons License CC-BY

10th International Conference on Interactive Theorem Proving (ITP 2019).

Editors: John Harrison, John O'Leary, and Andrew Tolmach; Article No. 13; pp. 13:1–13:19

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

the systems less direct. Moreover, we do not target automatic generation of executable code. We claim that our proof is direct, readable, elegant, and that it follows Tarjan's presentation. Note that a similar exercise but for a much more elementary proof (the irrationality of square root of 2) and using many more proof assistants (17) was presented in [35].

Examples of formal and informal proofs of algorithms about graphs can be found in [27, 33, 18, 28, 14, 20, 32, 22, 30, 29, 16, 8], among others. Some of them are part of a larger library, others focus on the treatment of pointers or on concurrent algorithms. In particular, only Lammich and Neumann [20] gave an alternative formal proof of Tarjan's algorithm within their framework for verifying graph algorithms in Isabelle/HOL.

We expose here the key parts of the proofs. The interested reader can access the details of the proofs and run them on the web [9, 10, 23]. In this paper, we recall the principles of the algorithm in Section 2; we describe the proofs in the three systems in Sections 3, 4, and 5 by emphasizing the differences induced by the logics which are used; we conclude in Sections 6 and 7 by commenting the developments and advantages of each proof system.

2 The algorithm

In a directed graph, two vertices x and y are strongly connected if there exists a path from x to y and a path from y to x . A strongly connected component (scc) is a maximal set of vertices where all pairs of vertices are strongly connected. The vertices reached by a depth-first search (DFS) traversal in a directed graph form a spanning forest. A fundamental property relates sccs and DFS traversal: each scc is a prefix of a single subtree in the spanning forest (see Figure 1c). Its root is called the base of the scc. Tarjan's algorithm [31] relies on the detection of these bases and collects the sccs in a pushdown stack. It performs a single DFS traversal of the graph assigning a serial number $num[x]$ to any vertex x in the order of the visit. It computes the following function for every vertex x :¹

$$LOWLINK(x) = \min\{num[y] \mid x \xRightarrow{*} z \hookrightarrow y \wedge x \text{ and } y \text{ are strongly connected}\}$$

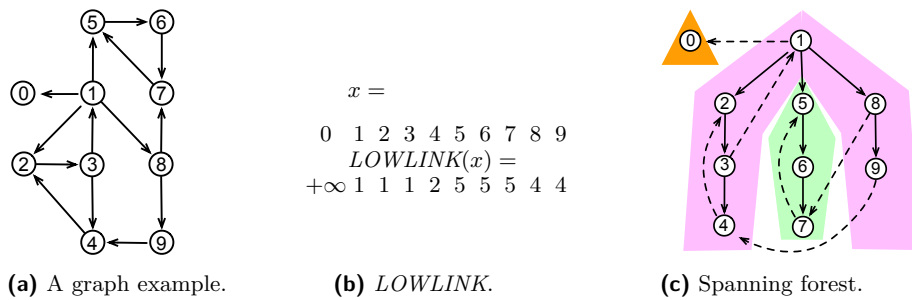
The relation $x \xRightarrow{*} z$ means that z is a son of x in the spanning forest, the relation $\xRightarrow{*}$ is its transitive and reflexive closure, and $z \hookrightarrow y$ means that there is a back edge from z to some node y of the spanning forest (a back edge is an edge of the graph which is not an edge in the spanning forest). In Figure 1c, $\xRightarrow{*}$ is drawn in thick lines and \hookrightarrow in dotted lines; in Figure 1b the table of the values of the *LOWLINK* function is shown. The minimum of the empty set is assumed to be $+\infty$ (this is a slight simplification w.r.t. the original algorithm).

The base x of an scc is found when $LOWLINK(x) \geq num[x]$, and the component is formed by the nodes of the subtree in the spanning forest rooted at x and pruned of the sccs already discovered in that subtree. As illustrated by vertices 8 or 9 in Figure 1c, notice that $LOWLINK(x)$ need not be the lowest serial number of a vertex accessible from x , nor of an ancestor of x in the spanning forest. The DFS traversal sets to $+\infty$ the serial numbers of vertices in already discovered sccs. This allows us to compute *LOWLINK* as:

$$LOWLINK(x) = \min\{num[y] \mid x \xRightarrow{*} z \hookrightarrow y\}$$

Our implementation of graphs uses an abstract type *vertex* for vertices, a constant *vertices* for the finite set of all vertices in the graph, and a *successors* function from vertices to their adjacency set. The algorithm maintains an environment e implemented as a record of type *env* with four fields: a stack *e.stack*, a set *e.sccs* of strongly connected components, a fresh serial number *e.sn*, and a function *e.num* from vertices to serial numbers.

¹ This definition relies on knowing whether nodes are strongly connected, which the algorithm is intended to discover. This apparent circularity will be broken below.



■ **Figure 1** The vertices are numbered and pushed onto the stack in the order of their visit by the recursive function *dfs1*. When the first component $\{0\}$ is discovered, vertex 0 is popped; similarly when the second component $\{5, 6, 7\}$ is found, its vertices are popped; finally all vertices are popped when the third component $\{1, 2, 3, 4, 8, 9\}$ is found. Notice that there is no back edge to a vertex with a number less than 5 when the second component is discovered. Similarly in the first component, there is no edge to a vertex with a number less than 0. In the third component, there is no edge to a vertex less than 1 since we have set the serial number of vertex 0 to $+\infty$ when 0 was popped.

```

type vertex
constant vertices: set vertex
function successors vertex : set vertex
type env = {stack: list vertex; sccs: set (set vertex); sn: int; num: map vertex int}

```

The DFS traversal is organized as two mutually recursive functions *dfs1* and *dfs*. The function *dfs1* visits a new vertex x and computes *LOWLINK*(x). Furthermore it adds a new scc when x is the base of a new scc. The function *dfs* takes as argument a set r of roots and an environment e . It calls *dfs1* on non-visited vertices in r and returns a pair consisting of an integer and the modified environment. The integer is the minimum of the values computed by *dfs1* on non-visited vertices in r and the serial numbers of already visited vertices in r . If the set of roots is empty, the returned integer is $+\infty$.

The main procedure *tarjan* initializes the environment with an empty stack, an empty set of sccs, the fresh serial number 0 and the constant function giving the number -1 to each vertex. The result is the set of components returned by the function *dfs* called on all vertices in the graph.

```

let rec dfs1 x e =
  let n0 = e.sn in
  let (n1, e1) = dfs (successors x) (add_stack_incr x e) in
  if n1 < n0 then (n1, e1) else
    let (s2, s3) = split x e1.stack in
    (+∞, {stack = s3; sccs = add (elements s2) e1.sccs;
          sn = e1.sn; num = set_infty s2 e1.num})
with dfs r e = if is_empty r then (+∞, e) else
  let x = choose r in let r' = remove x r in
  let (n1, e1) = if e.num[x] ≠ -1 then (e.num[x], e) else dfs1 x e in
  let (n2, e2) = dfs r' e1 in (min n1 n2, e2)
let tarjan () =
  let e = {stack = Nil; sccs = empty; sn = 0; num = const (-1)} in
  let (_, e') = dfs vertices e in e'.sccs

```

In the body of *dfs1*, the auxiliary function *add_stack_incr* updates the environment by pushing x on the stack, assigning it the current fresh serial number, and incrementing that number in view of future calls. The function *dfs1* performs a recursive call to *dfs* for the adjacent vertices of x as roots and the updated environment. If the returned integer value $n1$ is less than the number $n0$ assigned to x , the function simply returns $n1$ and the current

environment. Otherwise, the function declares that a new scc has been found, consisting of all vertices that are contained on top of x in the current stack. These vertices are popped from the stack, stored as a new set in $e.sccs$, and their serial numbers are all set to $+\infty$, ensuring that they do not interfere with future calculations of min values. The auxiliary functions *split* and *set_infty* are used to carry out these updates.

```

let add_stack_incr x e = let n = e.sn in
  {stack = Cons x e.stack; sccs = e.sccs; sn = n+1; num = e.num[x ← n]}
let rec set_infty s f = match s with Nil → f
  | Cons x s' → (set_infty s' f)[x ← +∞] end
let rec split x s = match s with Nil → (Nil, Nil)
  | Cons y s' → if x = y then (Cons x Nil, s')
    else let (s1', s2) = split x s' in (Cons y s1', s2) end

```

Figure 1 illustrates the behavior of the algorithm by an example. We presented the algorithm as a functional program, using data structures available in the Why3 standard library [4]. For lists we have the constructors *Nil* and *Cons*; the function *elements* returns the set of elements of a list. For finite sets, we have the empty set *empty*, and the functions *add* to add an element to a set, *remove* to remove an element from a set, *choose* to pick an arbitrary element² in a (non-empty) set, and *is_empty* to test for emptiness. We also use maps with functions *const* denoting the constant function, $[_]_[]$ to access the value of an element, and $[_]_[] \leftarrow []$ for creating a map obtained from an existing map by setting an element to a given value.

For a correspondence between our presentation and the imperative programs used in standard textbooks, the reader is referred to [8]. The present version can be directly translated into COQ or Isabelle functions, and it respects the linear running time of the algorithm for straightforward implementations of the elements that we choose to leave abstract in our presentation. Indeed, vertices could be represented by integers, $+\infty$ by the cardinal of *vertices*, finite sets by lists of integers and mappings by mutable arrays (see for instance [9]).

Thus for each environment e in the algorithm, the working stack $e.stack$ corresponds to a cut of the spanning forest where strongly connected components to its left are pruned and stored in $e.sccs$. In this stack, any vertex can reach any vertex higher in the stack. And if a vertex is a base of an scc, no back edge can reach some vertex lower than this base in the stack, otherwise that last vertex would be in the same scc with a strictly lower serial number.

Our proofs of the algorithm make these arguments formal. To maintain these invariants we will distinguish, as is common for DFS algorithms, three sets of vertices: white vertices are the non-visited ones, black vertices are those that are already fully visited, and gray vertices are those that are still being visited. Clearly, these sets are disjoint and white vertices can be considered as forming the complement in *vertices* of the union of the gray and black ones.

The previously mentioned invariant properties can now be expressed for vertices in the stack: no such vertex is white, any vertex can reach all vertices higher in the stack, any vertex can reach some gray vertex lower in the stack. Moreover, vertices in the stack respect the numbering order, i.e. a vertex x is lower than y in the stack if and only if the number assigned to x is strictly less than the number assigned to y .

3 The proof in Why3

The Why3 system comprises the programming language WhyML used in the previous section and a many-sorted first-order logic with inductive data types and inductive predicates to express the logical assertions. The system generates proof obligations w.r.t. the assertions,

² This is the only non-deterministic operation used in our development. We handle it by underspecification, effectively showing that any implementation is correct.

pre- and post-conditions and lemmas inserted in the WhyML program. The system is interfaced with off-the-shelf automatic provers and interactive proof assistants.

From the Why3 library, we use pre-defined theories for integer arithmetic, polymorphic lists, finite sets and mappings. There is also a small theory for paths in graphs. Here we define graphs, paths and sccs as follows.

```
axiom successors_vertices:  $\forall x$ . mem x vertices  $\rightarrow$  subset (successors x) vertices
predicate edge (x y: vertex) = mem x vertices  $\wedge$  mem y (successors x)
inductive path vertex (list vertex) vertex =
| Path_empty:  $\forall x$ : vertex. path x Nil x
| Path_cons:  $\forall x$  y z: vertex, l: list vertex.
  edge x y  $\rightarrow$  path y l z  $\rightarrow$  path x (Cons x l) z
predicate reachable (x y: vertex) =  $\exists l$ . path x l y
predicate in_same_scc (x y: vertex) = reachable x y  $\wedge$  reachable y x
predicate is_subsccl (s: set vertex) =  $\forall x$  y. mem x s  $\rightarrow$  mem y s  $\rightarrow$  in_same_scc x y
predicate is_scc (s: set vertex) = not is_empty s
 $\wedge$  is_subsccl s  $\wedge$  ( $\forall s'$ . subset s s'  $\rightarrow$  is_subsccl s'  $\rightarrow$  s == s')
```

The predicates *mem* and *subset* denote membership and the subset relation for finite sets.

We add two ghost fields in environments for the black and gray sets of vertices. These fields are used in the proofs but not used in the calculation of the sccs, which is checked by the type-checker of the language.

```
type env = {ghost black: set vertex; ghost gray: set vertex;
  stack: list vertex; sccs: set (set vertex); sn: int; num: map vertex int}
```

Defining a new function *add_black* turning a vertex from gray to black and redefining *add_stack_incr* for adding a new gray vertex with a fresh serial number to the current stack, the functions now become:

```
let add_black x e =
{black = add x e.black; gray = remove x e.gray;
 stack = e.stack; sccs = e.sccs; sn = e.sn; num = e.num}
let add_stack_incr x e =
let n = e.sn in
{black = e.black; gray = add x e.gray;
 stack = Cons x e.stack; sccs = e.sccs; sn = n+1; num = e.num[x  $\leftarrow$  n]}
let rec dfs1 x e =
let n0 = e.sn in
let (n1, e1) = dfs (successors x) (add_stack_incr x e) in
if n1 < n0 then (n1, add_black x e1) else
let (s2, s3) = split x e1.stack in
(+ $\infty$ , {black = add x e1.black; gray = e1.gray; stack = s3;
 sccs = add (elements s2) e1.sccs; sn = e1.sn; num = set_infty s2 e1.num})
with dfs r e = ... (* unmodified *)
let tarjan () =
let e = {black = empty; gray = empty;
 stack = Nil; sccs = empty; sn = 0; num = const (-1)} in
let (_, e') = dfs vertices e in e'.sccs
```

The main invariant (\mathcal{I}) of our program states that the environment is well-formed:

```
predicate wf_env (e: env) =
let {stack = s; black = b; gray = g} = e in
wf_color e  $\wedge$  wf_num e  $\wedge$  simplelist s  $\wedge$  no_black_to_white b g  $\wedge$ 
( $\forall x$  y. lmem x s  $\rightarrow$  lmem y s  $\rightarrow$  e.num[x]  $\leq$  e.num[y]  $\rightarrow$  reachable x y)  $\wedge$ 
( $\forall y$ . lmem y s  $\rightarrow$   $\exists x$ . mem x g  $\wedge$  e.num[x]  $\leq$  e.num[y]  $\wedge$  reachable y x)  $\wedge$ 
( $\forall cc$ . mem cc e.sccs  $\leftrightarrow$  subset cc b  $\wedge$  is_scc cc)
```

where *lmem* stands for membership in a list. The well-formedness property is the conjunction of seven clauses. The two first clauses express elementary conditions about the colored sets of vertices and the numbering function (see [9, 8] for a detailed description). The third clause

13:6 Formal Proofs of Tarjan's SCC Algorithm

states that there are no repetitions in the stack, and the fourth that there is no edge from a black vertex to a white vertex. The next two clauses formally express the property already stated above: any vertex in the stack reaches all higher vertices and any vertex in the stack can reach a lower gray vertex. The last clause states that the *sccs* field is the set of all sccs all of whose vertices are black.

Since at the end of the *tarjan* function, all vertices are black, the *sccs* field will contain exactly the set of all strongly connected components. We formally express this by adding a post-condition to the definition of the function.

```
let tarjan () = returns {r → ∀cc. mem cc r ↔ subset cc vertices ∧ is_scc cc}
  let e = {black = empty; gray = empty;
          stack = Nil; sccs = empty; sn = 0; num = const (-1)} in
  let (_, e') = dfs vertices e in assert {subset vertices e'.black};
    e'.sccs
```

Our functions *dfs1* and *dfs* modify the environment in a monotonic way. Namely, they augment the set of visited vertices (the black ones); they keep invariant the set of the ones currently under visit (the gray set); they increase the stack with new black vertices; they also discover new sccs and they keep invariant the serial numbers of vertices in the stack,

```
predicate subenv (e e': env) =
  subset e.black e'.black ∧ e.gray == e'.gray
  ∧ (∃s. e'.stack = s ++ e.stack ∧ subset (elements s) e'.black)
  ∧ subset e.sccs e'.sccs ∧ (∀x. lmem x e.stack → e.num[x] = e'.num[x])
```

Once these invariants are expressed, it remains to locate them in the program text and to add assertions which help to prove them. The pre-conditions of *dfs1* are quite natural: the vertex x must be a white vertex of the graph, and it must be reachable from all gray vertices. Moreover invariant (\mathcal{I}) must hold. The post-conditions of *dfs1* are the following. Firstly, (\mathcal{I}) and the monotonicity property *subenv* hold of the resulting environment. Second, vertex x is black at the end of *dfs1*. Finally we express properties of the integer value n returned by this function, which is indeed $LOWLINK(x)$, as noted previously. Formally, we give three properties for characterizing n . The returned value is never higher than the number of x in the final environment. Also, the returned value is either $+\infty$ or the number of a vertex in the stack reachable from x . Finally, if there is a (back) edge from a vertex y' in the new part of the stack to a vertex y in its old part, the returned value n must be lower or equal to the serial number of y .

```
let rec dfs1 x e =
  (* pre-condition *)
  requires {mem x vertices ∧ not mem x (union e.black e.gray)}
  requires {∀y. mem y e.gray → reachable y x}
  requires {wf_env e} (* I *)
  (* post-condition *)
  returns {(_, e') → wf_env e' ∧ subenv e e'}
  returns {(_, e') → mem x e'.black}
  returns {(n, e') → n ≤ e'.num[x]}
  returns {(n, e') → n = +∞ ∨ num_of_reachable_in_stack n x e'}
  returns {(n, e') → ∀y. xedge_to e'.stack e.stack y → n ≤ e'.num[y]}
```

The auxiliary predicates used above are formally defined in the following way.

```
predicate num_of_reachable_in_stack (n: int) (x: vertex)(e: env) =
  ∃y. lmem y e.stack ∧ n = e.num[y] ∧ reachable x y
predicate xedge_to (s1 s3: list vertex) (y: vertex) =
  (∃s2. s1 = s2 ++ s3 ∧ ∃y'. lmem y' s2 ∧ edge y' y) ∧ lmem y s3
```

Notice that the definition of *xedge_to* fits the definition of *LOWLINK* when the back edge ends at a vertex residing in the stack before the call of *dfs1*. The pre- and post-conditions for the function *dfs* are quite similar up to a generalization to sets of vertices considered as the roots of the algorithm (see [9]).

We now add seven assertions in the body of the *dfs1* function to help the automatic provers. In contrast, the function *dfs* needs no extra assertions in its body. In *dfs1*, when the number *n0* of *x* is strictly greater than the number *n1* resulting from the call to its successors, the first assertion states that *n1* cannot be $+\infty$; it helps in proving the next assertion. The second assertion states that a lower gray vertex is reachable from *x* and that thus the scc of *x* is not fully black at the end of *dfs1*. In that assertion the inequality $y \neq x$ is redundant, but helps showing the *sccs* constraint at the end of *dfs1*. The first four assertions in the “else” branch show that the vertices on top of *x* in the current stack form a strongly connected component and that *x* is the base of that scc. The final assertion helps proving that the coloring constraint is preserved at the end of *dfs1*.

```

let n0 = e.sn in
let (n1, e1) = dfs (successors x) (add_stack_incr x e) in
if n1 < n0 then begin
  assert {n1 ≠ +∞};
  assert {∃y. y ≠ x ∧ mem y e1.gray ∧ e1.num[y] < e1.num[x] ∧ in_same_scc x y};
  (n1, add_black x e1) end
else
  let (s2, s3) = split x e1.stack in
  assert {is_last x s2 ∧ s3 = e.stack ∧ subset (elements s2) (add x e1.black)};
  assert {is_subsc (elements s2)};
  assert {∀y. in_same_scc y x → lmem y s2};
  assert {is_scc (elements s2)};
  assert {inter e.gray (elements s2) == empty};
  (+∞, {black = add x e1.black; gray = e.gray; stack = s3;
        sccs = add (elements s2) e1.sccs; sn = e.sn; num = set_infty s2 e1.num})

```

The function *inter* denotes set intersection, and *is_last* is defined below.

```

predicate is_last (x: α) (s: list α) = ∃s'. s = s' ++ Cons x Nil

```

All proofs are discovered by the automatic provers except for two proofs carried out interactively in COQ. One is the proof of the black extension of the stack (from predicate *subenv* in the post-condition of *dfs1*) in case $n1 < n0$. The provers could not find a suitable witness for the existential quantifier, although the COQ proof is quite short. The second COQ proof is the fifth assertion in the body of *dfs1*, which asserts that any *y* in the scc of *x* belongs to *s2*. It is a maximality assertion which states that the set *elements s2* is a complete scc. The proof of that assertion is by contradiction. If *y* is not in *s2*, there must be an edge from *x'* in *s2* to some *y'* not in *s2* such that *x* reaches *x'* and *y'* reaches *y*. There are three cases, depending on the position of *y'*. Case 1 is when *y'* is in *sccs*: this is not possible since *x* would then be in *sccs* which contradicts *x* being gray. Case 2 is when *y'* is an element of *s3*: the serial number of *y'* is strictly less than the one of *x* which is *n0*. If $x' \neq x$, the back edge from *x'* to *y'* contradicts $n1 \geq n0$ (post-condition 5); if $x' = x$, then *y'* is a successor of *x* and again it contradicts $n1 \geq n0$ (post-condition 3). Case 3 is when *y'* is white, then $x' \neq x$ is impossible since *x'* is then black in *s2* and would be the origin of a black-to-white edge to *y'*; if $x' = x$, then *y'* is not white by post-condition 2 of *dfs*.

Some quantitative information about the Why3 proof is listed in Table 1. Alt-Ergo 2.3 and CVC4 1.5 proved the bulk of the proof obligations.³ The proof uses 49 lemmas that were all proved automatically, but with an interactive interface providing hints to apply inlining,

³ In addition to the results reported in the table, Spass was used to discharge one proof obligation.

13:8 Formal Proofs of Tarjan's SCC Algorithm

■ **Table 1** Performance results with provers in Why3-0.88.3 (in seconds, on a 3.3 GHz Intel Core i5 processor). Total time is 341.15 seconds. The two last columns contain the numbers of verification conditions and proof obligations. Notice that there may be several VCs per proof obligation.

provers	Alt-Ergo	CVC4	E-prover	Z3	#VC	#PO
49 lemmas	1.91	26.11	3.33		70	49
split	0.09	0.16			6	6
add_stack_incr	0.01				1	1
add_black	0.02				1	1
set_infty	0.03				1	1
dfs1	77.89	150.2	19.99	13.67	79	20
dfs	4.71	3.52		0.26	58	25
tarjan	0.85				15	5
total	85.51	179.99	23.32	13.93	231	108

splitting, or induction strategies. This includes 13 lemmas on sets, 16 on lists, 5 on lists without repetitions, 3 on paths, 5 on sccs and 7 very specialized lemmas directly involved in the proof obligations of the algorithm. The lemma *xpath_xedge* states a critical condition on paths, reducing a predicate on paths to a predicate on edges. In fact, most of the Why3 proof works on edges, which are handled more robustly by the automatic provers than paths. Another important lemma is *subsc_after_last_gray*, which shows that the elements of the stack on top of the last gray vertex form a subset of an scc. This means that a variant of the program with the *split* call before the if-statement would have a simpler proof, but its time complexity would be non-linear. The two COQ proofs are 9 and 81 lines long (the COQ files of 677 and 680 lines include preambles that are automatically generated during the translation from Why3 to COQ). The interested reader is referred to [9] where the full proof is available.

The proof explained so far only showed the partial correctness of the algorithm. But after adding two lemmas about union and difference for finite sets, termination is automatically proved by the following lexicographic ordering on the number of white vertices and roots.

```
let rec dfs1 x e = variant {cardinal (diff vertices (union e.black e.gray)), 0}
with dfs r e = variant {cardinal (diff vertices (union e.black e.gray)), 1, cardinal r}
```

4 The proof in Coq

COQ is based on type theory and the calculus of constructions, a higher order lambda-calculus, to express formulae and proofs. Some basic notions of graph theory are provided by the Mathematical Components Library [21]. Our formalization is parameterized by a finite type V for the vertices and the function *successors* such that *successors* x is the adjacency set of any vertex x . The boolean *gconnect* $x y$ indicates that a path connects the vertex x to the vertex y . The function *equivalence_partition* of the library creates a partition of a set with respect to an equivalence relation. The set *gsccs* of the sccs of a graph is defined as follows, based on *gconnect* and the set of all vertices $[set: V]$:

```
Definition gsymconnect x y := gconnect x y && gconnect y x.
Definition gsccs := equivalence_partition gsymconnect [set: V].
```

Components are represented as sets of sets $\{set \{set V\}\}$. Various operations on sets are available. In this proof, we use singletons $[set x]$, unions $S_1 \cup S_2$, differences $S_1 \setminus S_2$, complements $\sim: S$ and unions of all sets in a set of sets *cover* S .

COQ offers several mechanisms for combining properties (boolean conjunction, propositional conjunction, record, inductive family), all of which have their own idiosyncracies. In order to make the presentation more readable for a non-COQ expert, we write them all with the n -ary propositional conjunction $[\wedge P_1, P_2, \dots \& P_n]$. We refer to [10] for the actual code.

The COQ formalization differs from the one in Why3: it uses natural numbers only and does not mention colors (white, gray and black). In particular, the number ∞ is defined as the cardinality of V , vertices with $\infty + 1$ as serial number correspond to the white vertices of the previous section and the environment is defined as a record with only two fields, a set of *scs* and the mapping assigning serial numbers to vertices:

```
Record env := Env {escs : {set {set V}}; num: {ffun V → nat}}.
```

Given an environment e , the set of visited vertices is *visited* e (the vertices with serial number less or equal to ∞), the current fresh serial number is *sn* e (the cardinality of the set of visited vertices), and the stack is *stack* e (the list of elements x which satisfy $\text{num } e \ x < \text{sn } e$, sorted by increasing serial number).

Another difference with the Why3 algorithm is the definition of *dfs1* and *dfs* as two separate, rather than two mutually recursive functions. As in the Why3 program, *dfs1* takes a vertex x , and *dfs* a set of vertices *roots*, in addition to an environment e . In order to prepare for the combination of the two functions, they also take function parameters that will represent the recursive calls.

```
Definition dfs1 dfs x e :=
  let: (n1, e1) as res := dfs (successors x) (visit x e) in
  if n1 < sn e then res else (∞, store (stack e1 \ stack e) e1).
Definition dfs dfs1 dfs (roots : {set V}) e :=
  if [pick x in roots] isn't Some x then (∞, e) else
  let: (n1, e1) := if num e x ≤ ∞ then (num e x, e) else dfs1 x e in
  let: (n2, e2) := dfs (roots \ [set x]) e1 in (minn n1 n2, e2).
```

The expression *visit* $x \ e$ represents the environment where x gets the next serial number, and *store* produces an environment that contains an additional strongly connected component.

Then, the two functions are glued together in a recursive function *rec* where the parameter k controls the maximal recursive height.

```
Fixpoint rec k r e := if k is k'.+1 then dfs (dfs1 (rec k')) (rec k') r e else (∞, e).
```

If k is not zero (i.e. it is a successor of some k'), *rec* calls *dfs* taking care that its parameters can only use recursive calls to *rec* with a smaller recursive height, here k' . This ensures termination. A dummy value is returned in the case where k is zero. Finally, the top level *tarjan* calls *rec* with the proper initial arguments.

```
Definition tarjan := let: (_, e) := rec (∞ * (∞.+2)) V (Env ∅ [ffun ⇒ ∞.+1]) in escs e.
```

Initially, the roots are all the vertices (V) and the environment has no component and all vertices are not visited (their number is $\infty + 1$). As both *dfs* and *dfs1* cannot be applied more than the number of vertices, the value $\infty * (\infty + 2)$ encodes the lexicographic product of the two maximal heights. It gives *rec* enough fuel to never encounter the dummy value so *tarjan* correctly terminates the computation. This allows us to separate the proof of the termination from the algorithm itself, and this last statement is of course proved formally later as theorem *rec_terminates*.

The invariants of the COQ proof are usually shorter than in the Why3 proof since they do not mention colors. We first define well-formed environments and their valid extension:

13:10 Formal Proofs of Tarjan's SCC Algorithm

```

Definition wf_env e := [∧ escs e ⊆ gscs,
  ∀x, num e x < ∞ → num e x < sn e,
  ∀x, (num e x = ∞) = (x ∈ cover (escs e)) &
  ∀x y, num e x ≤ num e y < sn e → gconnect x y].
Definition subenv e1 e2 := [∧ escs e1 ⊆ escs e2,
  ∀x, num e1 x < ∞ → num e2 x = num e1 x &
  ∀x, num e2 x < sn e1 → num e1 x < sn e1].

```

Then we state that new visited vertices are the ones reachable by paths accessible from roots with non-visited vertices (i.e. by white paths in the colored setting). The function *nexts* such that *nexts D X* returns the set of vertices reachable from the set *X* by a path which only contains vertices in *D* except maybe the last one.

```

Definition outenv (roots : {set V}) (e e' : env) := [∧
  ∀x y, x ∈ stack e' \ stack e → y ∈ stack e' \ stack e → gconnect x y,
  ∀x, x ∈ stack e' \ stack e → ∃y, y ∈ stack e ∧ gconnect x y &
  visited e' = visited e ∪ nexts (λ: visited e) roots ].

```

The post-condition is the conjunction of these three properties and the characterization of the output rank:

```

Definition dfs_spec (ne' : nat * env) (roots : {set V}) e := let: (n, e') := ne' in
  [∧ n = \min_(x in nexts (λ: visited e) roots) inord (num e' x),
  wf_env e', subenv e e' & outenv roots e e'].

```

Here, the argument *ne'* is the result of a *dfs*. The output rank *n* is the minimum of the serial numbers of the vertices which can be reached from the roots through a path where all the vertices except maybe the last one were not already visited. Note that this characterization differs from the notion of *LOWLINK*, which requires that the last vertex was visited.

Finally, we express correctness as the implication between pre- and post-conditions:

```

Definition dfs_correct dfs (roots : {set V}) e := wf_env e →
  (∀x y, x ∈ stack e → y ∈ roots → gconnect x y) → dfs_spec (dfs roots e) roots e.
Definition dfs1_correct dfs1 x e := wf_env e → x ∉ visited e →
  (∀x y, x ∈ stack e → y ∈ [set x] → gconnect x y) → dfs_spec (dfs1 x e) [set x] e.

```

Although these invariants are expressed differently from the formulation in Why3, they reflect essentially the same ideas. Compared to a first version of the formal model that was more closely aligned with the Why3 representation, this more abstract version made it possible to reduce by approximately 50% the size of the Coq proofs. The two central theorems are:

```

Lemma dfsP dfs1 dfsrec (roots : {set V}) e : (∀x, x ∈ roots → dfs1_correct dfs1 x e) →
  (∀x, x ∈ roots → ∀e1, subenv e e1 → dfs_correct dfsrec (roots \ [set x]) e1) →
  dfs_correct (dfs dfs1 dfsrec) roots e.
Lemma dfs1P dfs x e : dfs_correct dfs (successors x) (visit x e) →
  dfs1_correct (dfs1 dfs) x e.

```

They state that *dfs* and *dfs1* are correct if their respective recursive calls are correct. The proof of the first lemma is straightforward since *dfs* simply iterates on a list. It mostly requires book-keeping between what is known and what needs to be proved. This is done in about 54 lines. The second one is more intricate and requires 124 lines. Gluing these two theorems together and proving termination gives us an extra 12 lines to prove, and the correctness of *tarjan* then follows directly in 19 lines of straightforward proof.

■ **Table 2** Distribution of the numbers of lines of the 43 proofs in the file *tarjan_nocolors*.

Number of lines	1	2	3	4	5	6	11	12	16	19	54	124
Number of proofs	19	7	5	2	1	2	2	1	1	1	1	1

Theorem `rec_terminates` `k (roots : {set V}) e :`
`k ≥ #|~: visited e| * (∞.+1) + #|roots| → dfs_correct (rec k) roots e.`
Theorem `tarjan_correct` : `tarjan = gscs`.

We now provide some quantitative information. The COQ contribution consists of two files. Module *extra_nocolors* defines the *bigmin* operator and some notions of graph theory that we intend to add to Mathematical Components. This file is 294 lines long. The main module is *tarjan_nocolors* and is 605 lines long. It is compiled in 12 seconds with a memory footprint of 800 Mb (3/4 of which are resident) on a Intel® i7 2.60GHz quad-core laptop running Linux. The proofs are performed in the SSREFLECT proof language [15] with very little automation. The proof script is mostly procedural, alternating book-keeping tactics (*move*) with transformational ones (mostly *rewrite* and *apply*), but often intermediate steps are explicitly declared with the *have* tactic. There are more than fifty of such intermediate steps in the 320 lines of proof of the file *tarjan_nocolors*. Table 2 gives the distribution of the numbers of lines of these proofs. Most of them are very short (26 are at most 2 lines long) and the only complicated proof is the one corresponding to the lemma *dfs1P*.

5 The proof in Isabelle/HOL

Isabelle/HOL [24] is the encoding of simply typed higher-order logic in the logical framework Isabelle [26]. Unlike Why3, it is not primarily intended as an environment for program verification and does not contain specific syntax for stating pre- and post-conditions or intermediate assertions in function definitions. Although logics and formalisms for program verification have been developed within Isabelle/HOL (e.g., [19]), we decided to express our formalization in plain Isabelle/HOL: our main objective is to compare a proof of a significant algorithm in different proof assistants. The constructions that we use are universally available and do not rely on any elaborate infrastructure that would make comparisons more difficult.

We start by introducing a *locale*, fixing parameters and assumptions for the remainder of the proof. Unlike in Why, where the *set* type constructor represents finite sets, we explicitly assume that the set of vertices is finite.

```
locale graph =
  fixes vertices :: ν set and successors :: ν ⇒ ν set
  assumes finite vertices and ∀v ∈ vertices. successors v ⊆ vertices
```

We introduce reachability using an inductive predicate definition, rather than via an explicit reference to paths as in Why3. Isabelle then generates appropriate induction theorems.

```
inductive reachable where
  reachable x x
  | [y ∈ successors x; reachable y z] ⇒ reachable x z
```

The definition of strongly connected components mirrors that used in Why3. The following lemma states that SCCs are disjoint; its one-line proof is found automatically using *Sledgehammer* [3], which heuristically selects suitable lemmas from the set of available facts (including Isabelle’s library), invokes several automatic provers, and finally reconstructs a proof that is checked by the Isabelle kernel.

13:12 Formal Proofs of Tarjan's SCC Algorithm

```

definition is_subsc where
  is_subsc S ≡ ∀x ∈ S. ∀y ∈ S. reachable x y
definition is_scc where
  is_scc S ≡ S ≠ {} ∧ is_subsc S ∧ (∀S'. S ⊆ S' ∧ is_subsc S' → S' = S)
lemma scc_partition:
  assumes is_scc S and is_scc S' and x ∈ S ∩ S'
  shows S = S'

```

Environments are represented by records with the same components as in Why3, and the definition of the well-formedness predicate is also essentially identical to Why3.⁴

```

record ν env =
  black :: ν set      gray :: ν set
  stack :: ν list    sccs :: ν set set  sn :: nat      num :: ν ⇒ int
definition wf_env where wf_env e ≡
  wf_color e ∧ wf_num e ∧ distinct (stack e) ∧ no_black_to_white e
  ∧ (∀x y. y ⋖ x in (stack e) → reachable x y)
  ∧ (∀y ∈ set (stack e). ∃g ∈ gray e. y ⋖ g in (stack e) ∧ reachable y g)
  ∧ sccs e = { C . C ⊆ black e ∧ is_scc C }

```

We now define the two mutually recursive functions *dfs1* and *dfs* that expect as arguments a vertex *x* and a set of vertices *roots*, as well as an environment.

```

function (domintros) dfs1 and dfs where
  dfs1 x e =
    (let (n1,e1) = dfs (successors x) (add_stack_incr x e) in
     if n1 < int (sn e) then (n1, add_black x e1)
     else (let (l,r) = split_list x (stack e1) in
           (+∞, (| black = insert x (black e1), gray = gray e,
                 stack = r, sn = sn e1, sccs = insert (set l) (sccs e1),
                 num = set_infty l (num e1) | ))) and
  dfs roots e =
    (if roots = {} then (+∞, e)
     else (let x = SOME x. x ∈ roots;
           res1 = (if num e x ≠ -1 then (num e x, e) else dfs1 x e);
           res2 = dfs (roots - {x}) (snd res1)
           in (min (fst res1) (fst res2), snd res2) ))

```

The **function** keyword introduces the definition of a recursive function. Isabelle checks that the definition is well-formed and generates appropriate simplification and induction theorems. Because HOL is a logic of total functions, two proof obligations are introduced: the first one requires the user to prove that the cases in the function definitions cover all type-correct arguments; this holds trivially for the above definitions. The second obligation requires exhibiting a well-founded ordering on the function parameters that ensures the termination of recursive function invocations, and Isabelle provides a number of heuristics that work in many cases. However, the functions defined above will in fact not terminate for arbitrary calls, in particular for environments that assign the serial number -1 to non-white vertices. The *domintros* attribute instructs Isabelle to consider these functions as “partial”. More precisely, it introduces an auxiliary predicate that represents the domains for which the functions are defined. This “domain condition” appears as a hypothesis in the simplification rules that mirror the function definitions. In particular, a function call can be replaced by the right-hand side of the definition only if the domain predicate holds. Isabelle also introduces (mutually inductive) rules for proving when the domain condition is known (or assumed) to hold. Our first objective is therefore to establish sufficient conditions that ensure

⁴ We use the infix operator \preceq to denote precedence in lists.

the termination of the two functions. Assuming the domain condition, we prove that the functions never decrease the set of colored vertices and that vertices are never explicitly assigned the number -1 by our functions. Denoting the union of gray and black vertices as *colored*, we show that the algorithm only colors vertices and never assigns them -1 . We then prove that the triples

```
(vertices - colored e, {x}, 1)
(vertices - colored e, roots, 2)
```

for the arguments of *dfs1* and *dfs*, respectively, decrease w.r.t. lexicographical ordering on finite subset inclusion and $<$ on natural numbers across recursive function calls, provided that *colored_num* holds when the function is called and that x is a white vertex. These conditions are therefore sufficient to ensure that the domain condition holds:⁵

```
theorem dfs1_dfs_termination:
  [|x ∈ vertices - colored e; colored_num e|] ⇒ dfs1_dfs_dom (Inl(x,e))
  [|roots ⊆ vertices; colored_num e|] ⇒ dfs1_dfs_dom (Inr(roots,e))
```

The proof of partial correctness follows the same ideas as the proof presented for Why3. We define the pre- and post-conditions of the two functions as predicates in Isabelle. For example, the two predicates for *dfs1* are defined as follows:

```
definition dfs1_pre where dfs1_pre e ≡
  wf_env e ∧ x ∈ vertices ∧ x ∉ colored e ∧ (∀g ∈ gray e. reachable g x)
definition dfs1_post where dfs1_post x e res ≡
  let n = fst res; e' = snd res
  in wf_env e' ∧ subenv e e' ∧ roots ⊆ colored e'
  ∧ (∀x ∈ roots. n ≤ num e' x)
  ∧ (n = +∞ ∨ (∃x ∈ roots. ∃y in set (stack e'). num e' y = n ∧ reachable x y))
```

We now prove the following theorems:

- The pre-condition of each function establishes the pre-condition of every recursive call appearing in the body of that function. For the second recursive call in the body of *dfs* we also assume the post-condition of the first recursive call.
- The pre-condition of each function, plus the post-conditions of each recursive call in the body of that function, establishes the post-condition of the function.

Combining these results, we establish partial correctness:

```
theorem dfs_partial_correct:
  [|dfs1_dfs_dom (Inl(x,e)); dfs1_pre x e|] ⇒ dfs1_post x e (dfs1 x e)
  [|dfs1_dfs_dom (Inr(roots,e)); dfs_pre roots e|] ⇒ dfs_post roots e (dfs roots e)
```

We now define the initial environment and the overall function. It is then trivial to show that the arguments to the call of *dfs* in the definition of *tarjan* satisfy the pre-condition of *dfs*. Putting together the theorems establishing termination and partial correctness, we obtain the desired total correctness results.

```
definition init_env where init_env ≡
  (| black = {}, gray = {}, stack = [], sccs = {}, sn = 0, num = λ_. -1 |)
definition tarjan where tarjan ≡
  sccs (snd (dfs vertices init_env))
theorem dfs_correct:
  dfs1_pre x e ⇒ dfs1_post x e (dfs1 x e)
  dfs_pre roots e ⇒ dfs_post roots e (dfs roots e)
theorem tarjan_correct:
  tarjan = { C . is_scc C ∧ C ⊆ vertices }
```

⁵ Observe that Isabelle introduces a single predicate *dfs1_dfs_dom* corresponding to the two mutually recursive functions whose domain is the disjoint sum of the domains of both functions.

■ **Table 3** Distribution of interactions in the Isabelle proofs.

$i = 1$	$i \leq 5$	$i \leq 10$	$i \leq 20$	$i \leq 30$	$i = 35$	$i = 43$	$i = 48$
28	8	4	1	2	1	1	1

The intermediate assertions that were inserted in the Why3 code guided the overall proof in Isabelle: they are established either as separate lemmas or as intermediate steps within the proofs of the above theorems. Similarly to the COQ proof, the overall induction proof was explicitly decomposed into individual lemmas as laid out above. In particular, whereas Why3 identifies the predicates that can be used from the function code and its annotation with pre- and post-conditions, these assertions appear explicitly in the intermediate lemmas used in the proof of theorem *dfs_partial_correct*. The induction rules that Isabelle generated from the function definitions were helpful for finding the appropriate decomposition of the overall correctness proof.

We extensively used *sledgehammer* in the development of these proofs for invoking automatic back-end provers, including the superposition provers E, Spass, and Vampire, and the SMT solvers CVC4 and Z3. Nevertheless, we found that in comparison to Why3, significantly more user interactions were necessary in order to guide the proof. On several occasions, the external back-end reported finding a proof, but the subsequent attempt to reconstruct a proof in Isabelle failed: *sledgehammer* mainly records the lemmas that are used by the automatic back-end and then invokes proof tools (such as *metis*) that can generate a detailed proof that can be certified by the Isabelle kernel, but these tools may not find a proof even when the original automatic prover succeeded. When automatic proof fails, the user has to decompose the proof into smaller steps. Although decomposition was often straightforward, a few steps were not so obvious and required designing a rather detailed proof strategy. Table 3 indicates the distribution of the number of interactions used for the proofs of the 46 lemmas the theory contains. These numbers cannot be compared directly to those shown in Table 2 for the COQ proof because an Isabelle interaction is typically much coarser-grained than a line in a COQ proof. As in the case of Why3 and COQ, the proofs of partial correctness of *dfs1* (split into two lemmas following the case distinction) required the most effort. It took about one person-month to carry out the case study, starting from an initial version of the Why3 proof. Processing the entire Isabelle theory on a laptop with a 2.7 GHz Intel® Core i5 (dual-core) processor and 8 GB of RAM takes 35 CPU seconds.

6 General comments about the proof

Our formal proofs refer to colors, finite sets, and the stack, although the informal correctness argument is about properties of strongly connected components in spanning trees. The algorithmician would explain the algorithm with spanning trees as in Tarjan’s article. It would be nice to extract a program from such a proof, but the program does not explicitly manipulate spanning trees, and the proof should be given in terms of variables and data that appear in the program.

A first version of the formal proof used *ranks* in the working stack and a flat representation of environments by adding extra arguments to functions for the black, gray, scc sets and the stack. The automatic provers of Why3 worked very well with this representation. But after remodelling the proof in COQ and Isabelle/HOL, it appeared to be cleaner to gather these extra arguments in records and have a single extra argument for environments. Also *ranks* disappeared in favor of the *num* function and the precedence relation, which are easier to understand. Why3’s automatic provers have more difficulties with the inlining of environments, but with a few hints they could still succeed.

Proving the correctness of Tarjan’s algorithm requires surprisingly few, and entirely elementary, concepts of finite graphs. With the exception of the use of the Mathematical Components library for COQ, we therefore did not use existing libraries formalizing advanced concepts of graph theory [12, 25].

When designing a formal representation of an algorithm, one has to decide at what level of abstraction the algorithm should be modeled. For example, the COQ formalization shows that one can represent Tarjan’s algorithm and proof using just serial numbers and the set of strongly connected components found so far, and that the stack used in the algorithm and the colors used in the proof can be reconstructed. The Why3 and Isabelle representations make these elements explicit: coloring of vertices is frequently used when reasoning about graph algorithms, and including the stack allows us to capture the linear time complexity of the algorithm.

There is always a tension between the concision of the proof, its clarity and its relation to the real program. Our presentation aimed at comparing different proof assistants, and we have allowed for a few redundancies while staying at the algorithmic level rather than capturing an implementation.

7 Conclusion

The formal proof expressed in this article was initially designed and implemented in Why3 [8] as the result of a long process, nearly a two-year half-time work with many attempts of proofs about various graph algorithms (depth first search, Kosaraju strong connectivity, bi-connectivity, articulation points, minimum spanning tree). Why3 has a clear separation between programs and the logic. It makes the correctness proof quite readable for a programmer. Also first-order logic is easy to understand. Moreover, one can prove partial correctness without caring about termination.

Another important feature of Why3 is its interface with various off-the-shelf theorem provers (mainly SMT provers). Thus the system benefits from the current technology in theorem provers. Clerical sub-goals can be delegated to these provers, which makes the overall proof shorter and easier to understand. Although the proof must be split in more elementary pieces, this has the benefit of improving its readability. Several hints about inlining or induction reasoning are still needed, and two COQ proofs were used. Technically, the automatic provers and the translations from the Why3 representation to their input languages are part of the trusted code base: proofs are not checked independently. The system records sessions and facilitates incremental proofs. However, the automatic provers are sometimes no longer able to handle a proof obligation after seemingly minor modifications to the formulation of the algorithm or the predicates, making the proof somewhat unstable.

The COQ and Isabelle proofs were inspired by the Why3 proof. Their development therefore required much less time although their text is longer. We do not know if the final proof would have been significantly different had the initial proof been developed in another system than Why3. The COQ proof uses SSREFLECT and the Mathematical Components library, which helps reduce the size of the proof compared to classical COQ. The proof also uses the bigops library and several other higher-order features which makes it more abstract.

In COQ, one could prove termination using well-foundedness [2, 5], but because of nested recursion the `Function` command fails, and both `Equations` and `Program Fixpoint` require the addition of an extra proof argument to the function. Instead, we define the functionals `dfs1` and `dfs` and recombine them in `rec` and `tarjan` by recursion on a natural number used as fuel. We prove partial correctness on functionals and postpone termination on `rec`.

■ **Table 4** Characteristics of the three formal systems for our case study.

	Why3	Coq	Isabelle/HOL
expressivity	-	+	+
readability	+	-	+
stability w.r.t changes	-	+	+
ease of use	-	-	-
automation	+	-	+
partial correctness vs. termination	+	-	-
trusted base	-	+	+
lines of automatic proof	395	0	314 ui
lines of manual proof	90	898	1690

Our COQ proof does not use significant automation.⁶ All details are explicitly expressed, but many of them were already present in the Mathematical Components library. Moreover, a proof certificate is produced and a functional program could in principle be extracted. The absence of automation makes the system very stable to use since the proof script is explicit, but it requires a higher degree of expertise from the user.

The Isabelle/HOL proof can be seen as a mid-point between the Why3 and COQ proofs. It uses higher order logic and the level of abstraction is close to the one of the COQ proof, although more readable in this case study. The proof makes use of Isabelle’s extensive support for automation. In particular, *sledgehammer* [3] was very useful for finding individual proof steps. It heuristically selects lemmas and facts available in the context and then calls automatic provers (SMT solvers and superposition-based provers for first-order logic). When one of these provers finds a proof, *sledgehammer* attempts to find a proof that can be certified by the Isabelle kernel, using various proof methods such as combinations of rewriting and first-order reasoning (*blast*, *fastforce* etc.), calls to the *metis* prover or reconstruction of SMT proofs through the *smt* proof method. Unlike in Why3, the automatic provers used to find the initial proof are not part of the trusted code base because ultimately the proof is checked by the kernel. The price to pay is that the degree of automation in Isabelle is still significantly lower compared to Why3. Adapting the proof to modified definitions was fast: the Isabelle/jEdit GUI eagerly processes the proof script and quickly indicates those steps that require attention.

The Isabelle proof also faces the termination problem to achieve general consistency. We chose to delay handling termination, using the *domintros* attribute. The proofs of termination and of partial correctness are independent; in particular, we obtain a weaker predicate ensuring termination than the one used for partial correctness. Although the basic principle of the termination proof is similar to the COQ proof and relies on considering functionals of which the recursive functions are fixpoints, its technical formulation based on an appropriate well-founded order is closer to informal arguments that programmers would give and avoids low-level reasoning about fuel.

One strong point of Isabelle/HOL is its nice L^AT_EX output and the flexibility of its parser, supporting mathematical symbols. Combined with the hierarchical Isar proof language [34], the proof is in principle understandable without actually running the system.

⁶ Hammers exist for COQ [11, 13] but unfortunately they currently perform badly when used in conjunction with the Mathematical Components library.

In the end, the three systems Why3, COQ, and Isabelle/HOL are mature, and each one has its own advantages w.r.t. readability, expressivity, stability, ease of use, automation, partial-correctness, code extraction, trusted base and length of proof (a subjective assessment appears in Table 4). Coming up with invariants that are both strong enough and understandable was by far the hardest part in this work. This effort requires creativity and understanding, although proof assistants provide some help: missing predicates can be discovered by understanding which parts of the proof fail. We think that formalizing the proof in all three systems was very rewarding and helped us better understand the state of the art in computer-aided deductive program verification. We would welcome further comparisons based on implementations of this quite challenging case study in other formal systems.⁷

Our work has not considered how the systems that we have used could have helped us generate an executable implementation of this algorithm, leave alone an efficient one, to imperative programs and concrete data structures. Formal refinement enables proceeding from the high-level correctness argument down to actually executable code, and there is support for verifying imperative programs in general-purpose proof assistants (e.g., [6, 7, 19]). However, the existing frameworks differ significantly, making comparisons quite difficult.

A final and totally different remark is about teaching of algorithms. Do we want students to formally prove algorithms, or to present algorithms with assertions, pre- and post-conditions, and make them prove these assertions informally as exercises? In both cases, we believe that our work could make a useful contribution.

References

- 1 João Alpuim and Wouter Swierstra. Embedding the refinement calculus in Coq. *Sci. Comput. Program.*, 164:37–48, 2018.
- 2 G. Barthe, J. Forest, D. Pichardie, and V. Rusu. Defining and reasoning about recursive functions: a practical tool for the Coq proof assistant. In *Functional and Logic Programming Systems (FLOPS'06)*, volume 3945 of *LNCS*, pages 114–129, Fuji Susono, Japan, April 2006.
- 3 Jasmin Christian Blanchette, Sascha Böhme, and Lawrence C. Paulson. Extending Sledgehammer with SMT Solvers. *J. Automated Reasoning*, 51(1):109–128, 2013.
- 4 François Bobot, Jean-Christophe Filliâtre, Claude Marché, Guillaume Melquiond, and Andrei Paskevich. *The Why3 platform, version 0.86.1*. LRI, CNRS & Univ. Paris-Sud & INRIA Saclay, version 0.86.1 edition, May 2015. Available at why3.lri.fr/download/manual-0.86.1.pdf.
- 5 Ana Bove, Alexander Krauss, and Matthieu Sozeau. Partiality and recursion in interactive theorem provers - an overview. *Mathematical Structures in Computer Science*, 26(1):38–88, 2016.
- 6 Arthur Charguéraud. Characteristic formulae for the verification of imperative programs. In *Proc. 16th ACM SIGPLAN Intl. Conf. Functional Programming*, pages 418–430, Tokyo, Japan, 2011. ACM.
- 7 Arthur Charguéraud. Higher-order Representation Predicates in Separation Logic. In *Proc. 5th ACM SIGPLAN Conf. Certified Programs and Proofs, CPP 2016*, pages 3–14, New York, NY, USA, 2016. ACM.
- 8 Ran Chen and Jean-Jacques Lévy. A Semi-automatic Proof of Strong connectivity. In *VSTTE 2017*, volume 10712 of *LNCS*, pages 49–65. Springer, 2017.
- 9 Ran Chen and Jean-Jacques Lévy. Full scripts of Tarjan SCC Why3 proof. Technical report, Iscas and Inria, 2017. jeanjacqueslevy.net/why3.

⁷ We have set up a Web page <http://www-sop.inria.fr/marelle/Tarjan/contributions.html> in order to collect formalizations.

- 10 Cyril Cohen and Laurent Théry. Full script of Tarjan SCC Coq/ssreflect proof, 2017. Available at <https://www-sop.inria.fr/marelle/Tarjan/>.
- 11 Łukasz Czajka and Cezary Kaliszzyk. Hammer for Coq: Automation for Dependent Type Theory. *J. Autom. Reasoning*, 61(1-4):423–453, 2018.
- 12 Christian Doczkal, Guillaume Combette, and Damien Pous. A Formal Proof of the Minor-Exclusion Property for Treewidth-Two Graphs. In *ITP 2018*, volume 10895 of *LNCS*, pages 178–195. Springer, 2018.
- 13 Burak Ekici, Alain Mebsout, Cesare Tinelli, Chantal Keller, Guy Katz, Andrew Reynolds, and Clark W. Barrett. SMTCoq: A plug-in for integrating SMT solvers into Coq. In *CAV (2)*, volume 10427 of *LNCS*, pages 126–133. Springer, 2017.
- 14 Jean-Christophe Filliâtre et al. The Why3 gallery of verified programs. Technical report, CNRS, Inria, U. Paris-Sud, 2015. toccata.lri.fr/gallery/why3.en.html.
- 15 Georges Gonthier and Assia Mahboubi. An introduction to small scale reflection in Coq. *J. Formalized Reasoning*, 3(2):95–152, 2010.
- 16 Aquinas Hobor and Jules Villard. The Ramifications of Sharing in Data Structures. In *Proc. 40th Ann. ACM SIGPLAN-SIGACT Symp. Principles of Programming Languages, POPL '13*, pages 523–536, New York, NY, USA, 2013. ACM.
- 17 Peter Lammich. Refinement for Monadic Programs. *Archive of Formal Proofs*, 2012. URL: https://www.isa-afp.org/entries/Refine_Monadic.shtml.
- 18 Peter Lammich. Verified Efficient Implementation of Gabow's Strongly Connected Component Algorithm. In *ITP 2015*, volume 8558 of *LNCS*, pages 325–340, Vienna, Austria, 2014. Springer.
- 19 Peter Lammich. Refinement to Imperative/HOL. *J. Automated Reasoning*, 2019.
- 20 Peter Lammich and René Neumann. A Framework for Verifying Depth-First Search Algorithms. In *Proc. 4th ACM SIGPLAN Conf. Certified Programs and Proofs, CPP '15*, pages 137–146, New York, NY, USA, 2015. ACM.
- 21 Assia Mahboubi and Enrico Tassi. *Mathematical Components*. Available at: <https://math-comp.github.io/mcb/>, 2016.
- 22 Farhad Mehta and Tobias Nipkow. Proving Pointer Programs in Higher-Order Logic. In *CADE*, 2003.
- 23 Stephan Merz. Formalization of Tarjan's Algorithm in Isabelle, 2018. Available at <https://members.loria.fr/SMerz/projects/tarjan/index.html>.
- 24 Tobias Nipkow, Lawrence Paulson, and Markus Wenzel. *Isabelle/HOL. A Proof Assistant for Higher-Order Logic*. Number 2283 in Lecture Notes in Computer Science. Springer Verlag, 2002.
- 25 Lars Noschinski. A Graph Library for Isabelle. *Mathematics in Computer Science*, 9(1):23–39, 2015.
- 26 Lawrence C. Paulson. *Isabelle: A Generic Theorem Prover*, volume 828 of *Lecture Notes in Computer Science*. Springer Verlag, 1994.
- 27 Christopher M. Poskitt and Detlef Plump. Hoare-Style Verification of Graph Programs. *Fundamenta Informaticae*, 118(1-2):135–175, 2012.
- 28 François Pottier. Depth-First Search and Strong Connectivity in Coq. In *Journées Francophones des Langages Applicatifs (JFLA 2015)*, January 2015.
- 29 Azalea Raad, Aquinas Hobor, Jules Villard, and Philippa Gardner. Verifying Concurrent Graph Algorithms. In *APLAS 2016*, volume 10017 of *LNCS*, pages 314–334, Hanoi, Vietnam, 2016. Springer.
- 30 Ilya Sergey, Aleksandar Nanevski, and Anindya Banerjee. Mechanized Verification of Fine-grained Concurrent Programs. In *Proc. 36th ACM SIGPLAN Conf. Programming Language Design and Implementation, PLDI '15*, pages 77–87, New York, NY, USA, 2015. ACM.
- 31 Robert Tarjan. Depth first search and linear graph algorithms. *SIAM Journal on Computing*, 1972.
- 32 Laurent Théry. Formally-proven Kosaraju's algorithm. Inria report, Hal-01095533, 2015.

- 33 Ingo Wengener. A Simplified Correctness Proof for a Well-Known Algorithm Computing Strongly Connected Components. *Information Processing Letters*, 83(1):17–19, 2002.
- 34 Markus Wenzel. Isar – A Generic Interpretative Approach to Readable Formal Proof Documents. In *TPHOLS'99*, volume 1690 of *LNCS*, pages 167–184, Nice, France, 1999. Springer.
- 35 Freek Wiedijk. *The Seventeen Provers of the World*, volume 3600 of *LNCS*. Springer, Berlin, Heidelberg, 2006.