

Formal methods and assisted proofs: application to strongly connected components algorithms

Département Informatique — Parcours Recherche
Tuteur : Stephan Merz

Vincent TRÉLAT

March 18, 2022

Contents

1	Introduction	3
1.1	Academic context	3
1.2	Formal methods	3
1.3	Isabelle (HOL)	3
1.4	Isabelle by example	3
2	Models and representation	4
2.1	Nodes	4
2.2	Graphs and their representation	4
3	Formalisation	5
3.1	Strongly connected components	5
3.1.1	Directed graphs	5
3.1.2	Examples	5
3.2	Order of traversal and backtracking edges	6
3.2.1	DFS and <i>num</i> value	6
3.2.2	Backtracking edges	7
3.3	Lowlink value	8
3.3.1	Definition	8
3.3.2	Example	8
4	Tarjan's algorithm	10
4.1	Description of Tarjan's algorithm	11
5	A sequential set-based algorithm	12
5.1	Formalisation	12
5.2	The algorithm	13
5.3	Informal proof	14
5.4	Formal proof	16
5.4.1	Environment setup	16
5.4.2	Implementation of the algorithm	17

1 Introduction

1.1 Academic context

This research work was carried out as part of my curriculum at the french [École des Mines de Nancy](#). All documents such as codes or source papers are available on a [GitHub repository](#).

1.2 Formal methods

Formal methods are a field of computer science related to mathematical logic and reasoning. The whole purpose of the discipline is to ensure by a logical proof that a given algorithm is not only correct on its domain of definition, but also to find – or define – that domain. Formal methods find applications in a variety of fields, both concrete, such as the railway industry or self-driving cars, and abstract, such as computational architecture.

Although a formal proof lies first on paper, the real formalisation starts when proofs are mechanised in a proof assistant.

1.3 Isabelle (HOL)

Isabelle is a generic proof assistant. It allows mathematical formulas to be expressed in a formal language and provides tools for proving those formulas in a logical calculus.

isabelle.in.tum.de

Isabelle is a really powerful low-level proof assistant coming with a higher order logic (HOL) proving environment making the proofs easily readable and comprehensible without adding any abstract overlay. The term "assistant" designates the fact that Isabelle has numerous tools allowing various automations in the proofs such as a theorem seeker or an automatic solver.

1.4 Isabelle by example

The following example is a good introduction to the use of Isabelle.

2 Models and representation

2.1 Nodes

Vertices of a graph can be represented as nodes. A node \mathcal{N} is a simple data structure composed of an index, a boolean value telling if it has already been visited and two integer values *num* and *lowlink* whose role will be explained later. In the following, the aforesaid attributes will be referred to through the following notation:

Let \mathcal{N} be a node. The attributes of \mathcal{N} can be accessed via `N.index`, `N.visited`, `N.num` and `N.lowlink`. This notation will be applied to any object that lends itself to it.

2.2 Graphs and their representation

A graph \mathcal{G} is the data $(\mathcal{V}, \mathcal{E})$ where:

- \mathcal{V} is a set of vertices
- $\mathcal{E} \subseteq \{(x, y) \in V^2\}$ is a set of edges¹.

Vertices will often be called nodes and edges will be represented through adjacency lists for each node.

Let us give an example. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be the graph represented on figure 1. Thus, $\mathcal{V} = \{0, 1, 2, 3, 4, 5, 6, 7\}$ and

$$\begin{aligned} \mathcal{E} = \{ & (0, 0), (0, 1), (0, 2), (0, 3), \\ & (1, 4), (1, 7), \\ & (3, 0), (3, 1), (3, 2), (3, 5), \\ & (4, 3), (4, 6), \\ & (5, 6), \\ & (6, 3), \\ & (7, 6) \} \end{aligned}$$

This representation being somewhat long, adjacency lists can be used instead and therefore it gives:

$$\mathbf{G.adjacency} = [[0, 1, 2, 3], [4, 7], [0, 1, 2, 5], [3, 6], [6], [3], [6]]$$

Thus, for all $i \in \{0, \dots, 7\}$, `G.adjacency[i]` is the list of nodes to which node i is connected – *i.e.* there is an directed edge from node i to every node of `G.adjacency[i]`.

¹Note the use of the couple (x, y) and not the pair $\{x, y\}$ that makes the graph directed.

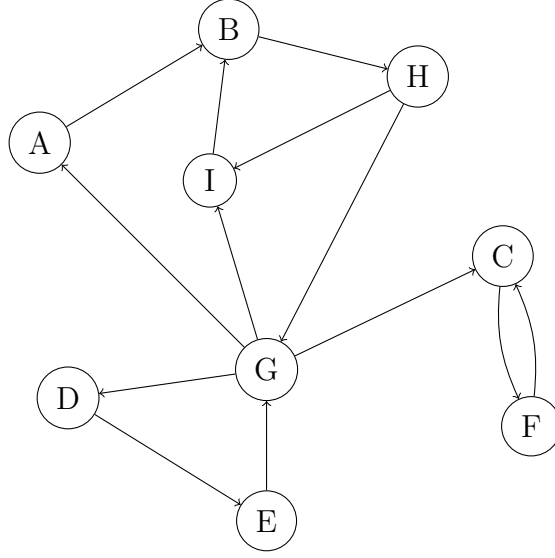


Figure 1: Representation of \mathcal{G}

3 Formalisation

3.1 Strongly connected components

3.1.1 Directed graphs

Definition 1. For two vertices x and y of \mathcal{V} , the relation "has an edge to" is noted " \Rightarrow " such that

$$(x, y) \in \mathcal{E} \iff x \Rightarrow y$$

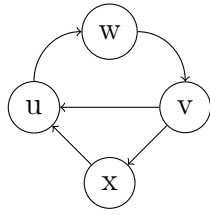
By extension, if there is a path from x to y with more than one edge, the same notation is kept for the sake of simplicity. The reflexive and transitive closure of the relation \Rightarrow is noted \Rightarrow^* .

Definition 2. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be an directed graph. $\mathcal{C} \subseteq \mathcal{V}$ is a strongly connected component of \mathcal{G} if:

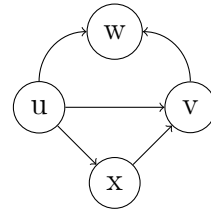
$$\forall x, y \in \mathcal{C}, (x \Rightarrow y) \wedge (y \Rightarrow x)$$

i.e. there is a path between every x and y in \mathcal{C} .

3.1.2 Examples



(a) Strongly connected graph



(b) Not strongly connected graph

Figure 2: Basic example of what is a small SCC

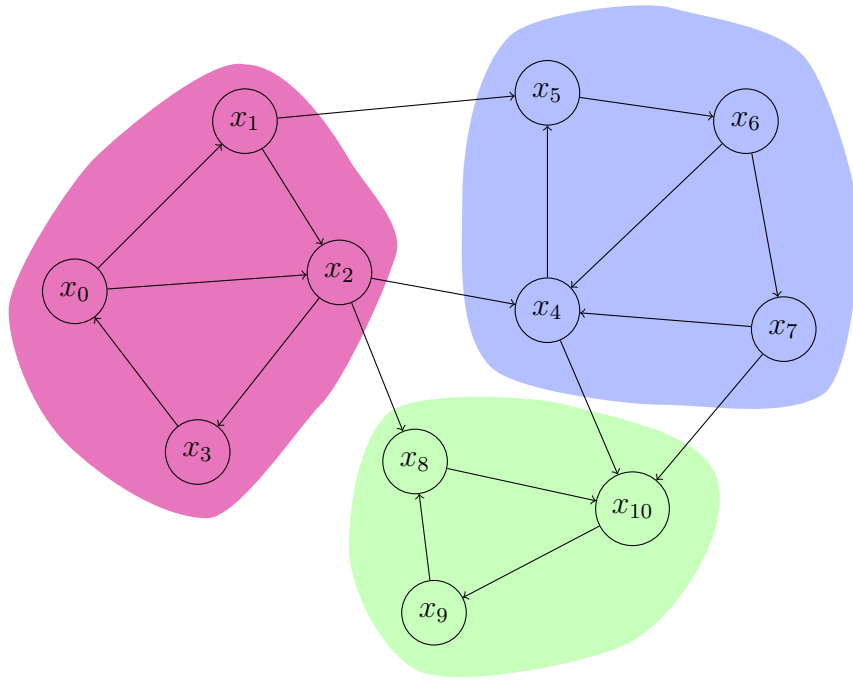


Figure 3: Example of a graph where each colored set of node is a maximal set of SCC

3.2 Order of traversal and backtracking edges

3.2.1 DFS and *num* value

Tarjan's SCC algorithm basically lies on a depth-first search. The figure 5 shows an example of a DFS traversal on a simple directed graph.

The previously mentioned figure also displays in red the **num** value which represents the order in which the nodes are visited in the graph during the DFS.

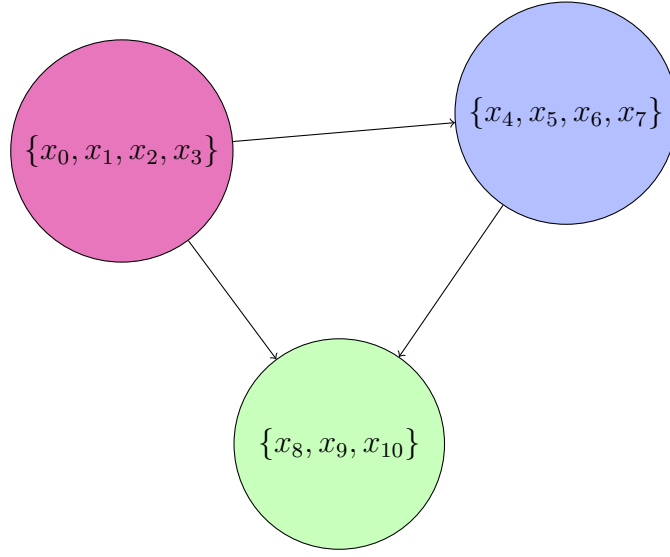


Figure 4: Reduced visualization of the graph represented if figure 3

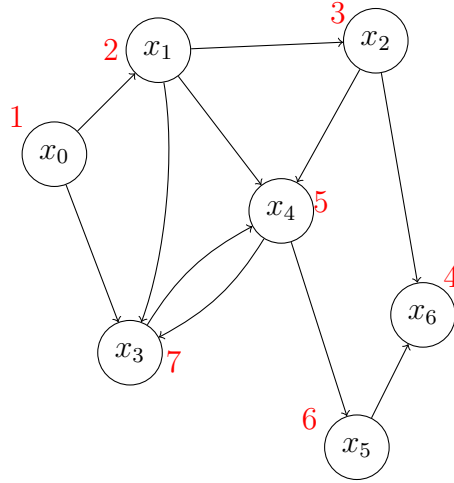


Figure 5: Example of a DFS

3.2.2 Backtracking edges

Definition 3. Given a graph \mathcal{G} and an order of traversal in this graph, *i.e.* each node of \mathcal{G} has a unique value $\text{num} \in [0, |\mathcal{V}|]$ and two nodes u and v , there is a backtracking edge from v to u if:

$$\begin{cases} u.\text{num} < v.\text{num} \\ v \Rightarrow u, \text{ i.e. } (v, u) \in \mathcal{E} \end{cases}$$

In this case, the backtracking edge from v to u is represented by $v \hookrightarrow u$.

3.3 Lowlink value

3.3.1 Definition

Informally, the `lowlink` value of a node represents the `num` value of the attachment node of their SCC, *i.e.* the `num` value of the entrance node in the corresponding SCC.

A more formal definition would be the following:

Definition 4. Let u be a node.

$$u.\text{lowlink} = \min\{w.\text{num} \mid \exists v \in \mathcal{V}, u \Rightarrow v \hookrightarrow w\}$$

3.3.2 Example

Let \mathcal{G} be the graph given in fig. 6. The order of traversal of the graph is given by the value *num* for each node of \mathcal{G} . The *lowlink* value is also displayed.

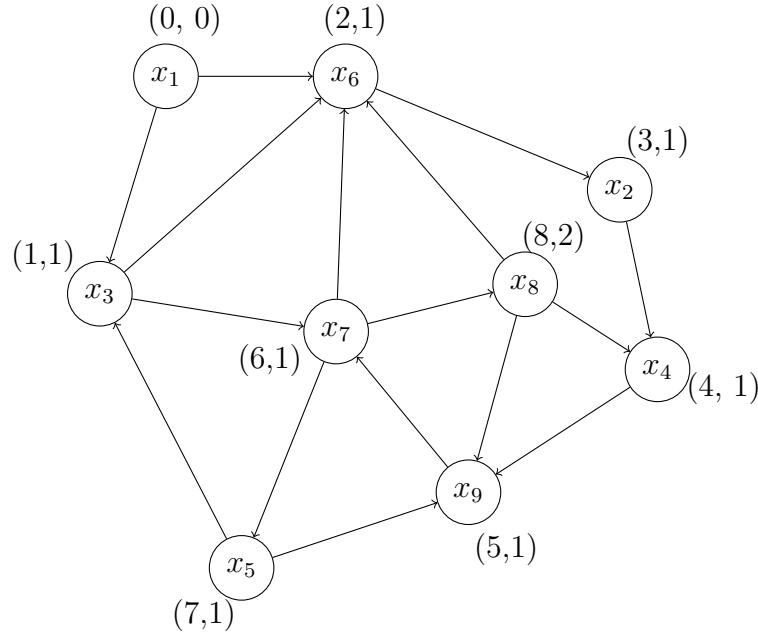


Figure 6: A DFS was performed through \mathcal{G} from x_1 and next to each node is represented the couple of value $(num, lowlink)$

Now, backtracking edges can be highlighted w.r.t. the order of traversal. In fig. 7, they are represented as red dashed arrows.

Knowing the backtracking edges, all *lowlink* values can be computed².

Let us take x_8 as an example: its *lowlink* is equal to 2, which actually means that x_6 ³ is its anchor – or attachment node – in their SCC, namely the green one. Indeed, x_8 is alone in its equivalence class⁴, and from all nodes linked by one of the backtracking edges of x_8 , x_6 has the minimum value *num*. Likewise, x_3 is its own attachment node in the green SCC since it is the first node visited when performing the DFS.

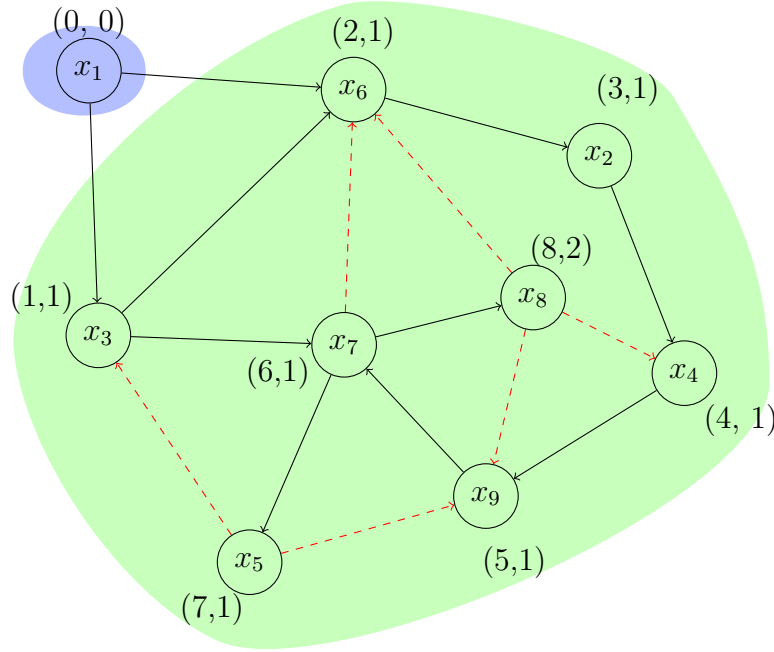


Figure 7: Same graph as in fig. 6 whose backtracking edges have been represented with red dashed arrows and SCCs have been highlighted

Then, SCCs can be easily found, namely $\{x_1\}$ and $\{x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9\}$, as shown in fig. 7.

²In fact, they are refreshed during the DFS.

³Because $x_6.\text{num} = 2$

⁴for the relation \Rightarrow

4 Tarjan's algorithm

Tarjan's algorithm is an efficient on-the-fly SCC computing algorithm [1]. It basically performs a DFS while updating the *num* and *lowlink* values. All nodes are stored in a stack during the traversal until a backtracking edge is found. In this case, the *lowlinks* are computed and all nodes are unstacked and saved in a SCC until a node verifying the equality between its *num* and its *lowlink* – which has to occur – is found. Then, the DFS goes on. The whole process is written in the following algorithm 1.

It can be shown that every node and edge are visited only once so the algorithm can achieve a linear complexity, *i.e.* $\mathcal{O}(|\mathcal{V}| + |\mathcal{E}|)$.

4.1 Description of Tarjan's algorithm

Algorithm 1: Tarjan's algorithm

Data: A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$

Result: A partition SCCs of \mathcal{V} where each element of SCCs is a SCC of \mathcal{G}

```
1 Initialize an empty stack R;
2 Initialize an empty set SCCs;
3 Let num := 0;
4 forall  $v \in \mathcal{V}$  do
5   if  $v.num$  is undefined then
6     SCC(v);
7   end
8 end
9 function  $SCC: v \in \mathcal{V} \rightarrow None$ 
10  v.num = num;
11  v.lowlink = num;
12  increment num;
13  Push  $v$  in R;
14  v.onStack = true;
15  forall  $w \in POST(v)$  do
16    if  $w.num$  is undefined then
17      v.lowlink = min(v.lowlink, w.lowlink);
18    end
19    else if  $w.onStack$  then
20      v.lowlink = min(v.lowlink, w.num);
21    end
22  end
23  if  $v.lowlink = v.num$  then
24    Initialize an empty set currentSCC;
25    repeat
26      Let  $w := R.pop()$ ;
27      w.onStack = false;
28      currentSCC = currentSCC  $\cup \{w\}$ ;
29    until  $v.num \neq w.num$ ;
30    SCCs = SCCs  $\cup$  currentSCC;
31  end
```

5 A sequential set-based algorithm

5.1 Formalisation

Definition 5 (SCC mapping). In the following algorithm, the SCCs are progressively tracked in a collection of disjoint sets through a map $\mathcal{S} : \mathcal{V} \rightarrow \mathcal{P}(\mathcal{V})$, where $\mathcal{P}(\mathcal{V})$ is the powerset of \mathcal{V} , s.t. the following invariant is maintained:

$$\forall v, w \in \mathcal{V}, w \in \mathcal{S}(v) \iff \mathcal{S}(v) = \mathcal{S}(w) \quad (1)$$

Remark 1. In particular, $\forall v \in \mathcal{V}, v \in \mathcal{S}(v)$.

Definition 6 (SCC union). Let UNITE be the function taking as parameters a map \mathcal{S} as defined previously and two vertices u and v of \mathcal{V} such that $\text{UNITE}(\mathcal{S}, u, v)$ merges the two mapped sets $\mathcal{S}(u)$ and $\mathcal{S}(v)$ and maintains the invariant (1) by updating the function \mathcal{S} .

Let us give an example:

Let $\mathcal{V} = \{u, v, w\}$ such that there is the following mapping: $\mathcal{S}(u) = \{u\}$ and $\mathcal{S}(v) = \mathcal{S}(w) = \{v, w\}$.

Then, $\text{UNITE}(\mathcal{S}, u, v) = \mathcal{S}(u) = \mathcal{S}(v) = \mathcal{S}(w) = \{u, v, w\}$.

Definition 7 (Successors set for a node). Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and $v \in \mathcal{V}$. The set of successors of v in \mathcal{G} is $\text{POST}(v)$ such that:

$$\forall w \in \text{POST}(v), (v, w) \in \mathcal{E}$$

5.2 The algorithm

See [3] for the original paper.

Algorithm 2: Sequential set-based SCC algorithm

Data: A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, a starting node v_0 ;
Result: A partition SCCs of \mathcal{V} where each element of SCCs is a maximal set of strongly connected components of \mathcal{G} ;

```

1 Initialize an empty set EXPLORED;
2 Initialize an empty set VISITED;
3 Initialize an empty stack R;
4 setBased( $v_0$ );
5 function setBased:  $v \in \mathcal{V} \rightarrow \text{None}$ 
6   VISITED := VISITED  $\cup \{v\}$ ;
7   R.push( $v$ );
8   foreach  $w \in \text{POST}(v)$  do
9     if  $w \in \text{EXPLORED}$  then
10       | continue;
11     end
12     else if  $w \notin \text{VISITED}$  then
13       | setBased( $w$ );
14     end
15     else
16       | while  $\mathcal{S}(v) \neq \mathcal{S}(w)$  do
17         |    $r := \text{R.pop}()$ ;
18         |   UNITE( $\mathcal{S}, r, \text{R.top}()$ );
19         | end
20       | end
21   end
22   if  $v = \text{R.top}()$  then
23     | report SCC  $\mathcal{S}(v)$ ;
24     | EXPLORED := EXPLORED  $\cup \mathcal{S}(v)$ ;
25     | R.pop();
26   end

```

5.3 Informal proof

Note that this proof is said informal only because it is not formally automated. Both logical and mathematical arguments developed below are absolutely relevant.

Lemma 1. (First invariant)

$$\forall x, y \in \mathbf{R}, x \neq y \implies \mathcal{S}(x) \cap \mathcal{S}(y) = \emptyset$$

Note the misuse of the set notation $x, y \in \mathbf{R}$ which just means that x and y are in the stack \mathbf{R} .

Proof. Let $x \in \mathcal{V}$ be the following node to be visited during the execution of the algorithm 2: x is pushed in \mathbf{R} . Let $y \in \text{POST}(x)$. There are two cases:

- y has not been visited yet, *i.e.* $y \notin \text{VISITED}$. Thus, a DFS-like traversal is performed from y , so y is pushed in \mathbf{R} and $\mathcal{S}(y) = \{y\}$ because y is alone in its equivalence class for the moment since it has not been visited yet. Therefore, $\mathcal{S}(x) \cap \mathcal{S}(y) = \emptyset$.
- y has already been visited, *i.e.* $y \in \text{VISITED}$. Then, y was already pushed in \mathbf{R} before x . Let $(x_i)_{1 \leq i \leq n}$ be the first nodes of the stack s.t. $x_0 = x$ and $x_n = y$.

In order to avoid writing $\mathbf{R} = [\dots, y, \dots, x]$, let us define \tilde{R} the stack containing the first n nodes in \mathbf{R} , s.t. $\tilde{R} = [y, \dots, x] = [x_n, \dots, x_0]$.

Let us consider the worst case, *i.e.* when

$$\forall 1 \leq i \leq n, \mathcal{S}(x_i) = \{x_i\}$$

So, the while loop has to go down to y because all partial SCCs are disjoint. As the length of the stack \mathbf{R} is bounded by $|\mathcal{V}|$, the algorithm terminates. x_0 is first unstacked and both $\mathcal{S}(x_0)$ and $\mathcal{S}(\mathbf{R}.\text{top}()) = \mathcal{S}(x_1)$ are then united. The current state of \mathcal{S} and \tilde{R} is:

$$\begin{cases} \mathcal{S} = \{\{x_0, x_1\}, \{x_2\}, \dots, \{x_n\}, \dots\} \\ \tilde{R} = [x_n, \dots, x_1] \end{cases}$$

Then, x_1 is unstacked and $\mathcal{S}(x_1)$ and $\mathcal{S}(x_2)$ are then united, so that:

$$\begin{cases} \mathcal{S} = \{\{x_0, x_1, x_2\}, \{x_3\}, \dots, \{x_n\}, \dots\} \\ \tilde{R} = [x_n, \dots, x_2] \end{cases}$$

Finally (by induction), $\mathcal{S} = \{x_0, \dots, x_n\}$ and $\tilde{R} = [y]$, *i.e.* $\mathcal{S}(x) = \mathcal{S}(y)$. It is important to notice that $x = x_0, x_1, \dots, x_{n-1}$ are no longer in the stack, so this operation kept the invariant true.

Lemma 2.

$$\biguplus_{v \in \mathbf{R}} \mathcal{S}(v) = \text{LIVE} := \text{VISITED} \setminus \text{EXPLORED}$$

Proof. The disjointness of all on-stack partial SCCs is given by lemma 1. Nodes from $\text{VISITED} \setminus \text{EXPLORED}$ are in \mathbf{R} because they are being processed. So, $\text{LIVE} \subseteq \mathbf{R}$.

By L.6-7 of algorithm 2, $\text{VISITED} \subseteq \mathbf{R}$.

L.9-10 ensure that no explored node is pushed in \mathbf{R} .

L.24-25 keep the invariant by unstacking explored nodes from \mathbf{R} , so $\mathbf{R} \cap \text{EXPLORED} = \emptyset$. Thus, $\mathbf{R} = \text{VISITED} \setminus \text{EXPLORED} = \text{LIVE}$. ■

Corollary 2.1.

$$\forall v \in \text{LIVE}, \exists! r \in \mathbf{R} \cap \mathcal{S}(v), \mathcal{S}(v) = \mathcal{S}(r)$$

Proof. Let $v \in \text{LIVE} = \biguplus_{v \in \mathbf{R}} \mathcal{S}(v)$. v is in a unique partial SCC $\mathcal{S} := \mathcal{S}(v)$. Because of lemma 1, there cannot exist $x \neq y \in \mathbf{R}$ s.t. $\mathcal{S}(x) = \mathcal{S}(y) = \mathcal{S}$. Thus, there exists a unique $x \in \mathbf{R}$ s.t. $\mathcal{S}(x) = \mathcal{S}$ (and $x \in \mathbf{R} \cap \mathcal{S}$). ■

Corollary 2.2.

$$\forall v \in \mathcal{V}, \forall w \in \text{POST}(v), w \in \text{LIVE} \implies \exists w' \in \mathbf{R}, \mathcal{S}(w') = \mathcal{S}(w)$$

Proof. Holds because of corollary 2.1. ■

Remark 2. In the algorithm 2, this property is held by L.16-18. These lines also illustrate how the algorithm "reads" the SCCs. Corollary 2.2 shows that when the mapped representatives of the top two nodes of \mathbf{R} are united (until $\mathcal{S}(w') = \mathcal{S}(v) = \mathcal{S}(w)$ since w' has a path to v), then all united components are in the same SCC.

Remark 3. Because \mathbf{R} only contains exactly one representative for each partial SCC (corollary 2.1), after each step of the main loop – *i.e.* the DFS – every partial SCC is actually maximal in the current set of visited nodes.

Theorem 1. The sequential algorithm 2 is correct, *i.e.* it returns a set of maximal SCCs.

Proof. Holds by remark 3. ■

5.4 Formal proof

Since the informal proof seems to be consistant, the formal – automated – proof can be written in Isabelle (HOL) based on the basis of the reasoning developped above.

5.4.1 Environment setup

The first definitions should be the different structures used in the algorithm. In particular, a record containing all the sets needed and described in the pseudo-code of algorithm 2 :

```
record 'v env =  
  S :: "'v  $\Rightarrow$  'v set"  
  explored :: "'v set"  
  visited :: "'v set"  
  sccs :: "'v set set"  
  stack :: "'v list"
```

The following lines define a graph structure and some useful natural relations :

```
locale graph =  
  fixes vertices :: "'v set" and successors :: "'v  $\Rightarrow$  'v set"  
  assumes vfin: "finite vertices"  
  and sclosed: " $\forall x \in \text{vertices}. \text{successors } x \subseteq \text{vertices}$ "  
  
  abbreviation edge where  
    "edge x y  $\equiv$  y  $\in$  successors x"  
  
  inductive reachable where  
    reachable_refl[iff]: "reachable x x"  
  | reachable_succ[elim]: "[edge x y; reachable y z]  $\implies$  reachable x z"
```

Those two relations are **edge**, which simply translates the property for two nodes of being linked by an edge and **reachable**, which is the binary relation \Rightarrow^* defined in section 3.1.1.

In order to be able to use those relations in the proofs later, it is essential to prove a list of lemmas, namely all the different natural properties that Isabelle cannot deduce⁵ from nothing⁶. For instance, the following lemmas are essential.

⁵That is an abuse of language. The idea is for example that for the moment, there is no formal link between **edge** and **reachable**. The goal is to formalize it so Isabelle is logically able to both use and simplify some results in the proofs.

⁶There is actually a theorem fetcher that is particularly useful to find a basic set of lemmas.

```
lemma reachable_edge: "edge x y  $\implies$  reachable x y"
  by auto
```

Mathematical writing: $\forall x, \forall y, x \Rightarrow y \implies x \Rightarrow^* y$

```
lemma succ_reachable:
  assumes "reachable x y" and "edge y z"
  shows "reachable x z"
  using assms by induct auto
```

Mathematical writing: $\forall x, \forall y, \forall z, (x \Rightarrow^* y \wedge y \Rightarrow z) \implies x \Rightarrow^* z$

```
lemma reachable_trans:
  assumes y: "reachable x y" and z: "reachable y z"
  shows "reachable x z"
  using assms by induct auto
```

Mathematical writing: $\forall x, \forall y, \forall z, (x \Rightarrow^* y \wedge y \Rightarrow^* z) \implies x \Rightarrow^* z$

As the formal proofs will eventually deal with strongly connected components, it is also essential to formally define SCCs. For the purpose of the proof, the property of being a SCC is called `sub_scc` and being a *maximal* SCC is called `is_scc` :

```
definition is_subscs where
```

```
"is_subscs S  $\equiv$   $\forall x \in S. \forall y \in S. \text{reachable } x y$ "
```

Mathematical writing: A set S is a SCC if $\forall x \in S, \forall y \in S, x \Rightarrow^* y$

```
definition is_scc where
```

```
"is_scc S  $\equiv$  S  $\neq$  {}  $\wedge$  is_subscs S
 $\wedge$  ( $\forall S'. S \subseteq S' \wedge \text{is\_subscs } S' \longrightarrow S' = S$ )"
```

Mathematical writing: A non-empty SCC S is maximal if for all SCC S' , $S \subseteq S' \implies S' = S$

Once again, there are some lemmas to prove, such as telling Isabelle when an element can be added to a SCC, or that two vertices that are reachable from each other are in the same SCC, or that two SCCs having a common element are identical.

5.4.2 Implementation of the algorithm

Now that the environment is set up, the actual algorithm – seen as a function – can be implemented.

Since Isabelle does not support loops, the implementation will be split into two mutually recursive functions. The main function is called `dfs` and takes its name after the Depth First Search algorithm because the algorithm 2 roughly consists in a deep traversal of a graph. The second function is called `dfss` and represents the *while* loop of the algorithm

2. The two functions are mutually recursive because they recursively call each other. In particular, `dfss` will call either itself or `dfs`, depending on the case.

```
function dfs :: "'v ⇒ 'v env ⇒ 'v env" and
  dfss :: "'v ⇒ 'v set ⇒ 'v env ⇒ 'v env" where
"dfs v e =
  (let e1 = e(|visited := visited e ∪ {v}, stack := (v # stack e)|);
   e' = dfss v (successors v) e1
   in if v = hd(stack e')
      then e'(|sccs:=sccs e' ∪ S e' v, explored:=explored e' ∪ (S e' v),
stack:=tl(stack e')|)
      else e')"
| "dfss v vs e =
  (if vs = {} then e
   else (let w = SOME x. x ∈ vs
        in (let e' = (if w ∈ explored e then e
                      else if w ∉ visited e then dfs w e
                      else unite v w e)
            in dfss v (v - {w}) e'))))"
by pat_completeness (force+)
```

The two last keywords require explanations as well : `pat_completeness` stands for *pattern completeness* and ensures that there is no missing patterns. The keyword `force` is used⁷ to help Isabelle know – by proving it – that both `dfs` and `dfss` are actually functions and that those functions are well defined with respect to the usual logical and mathematical meaning.

⁷Here, `auto` cannot terminate because of the mutual recursion.

References

- [1] R. Chen, C. Cohen, J.-J. Lévy, S. Merz, L. Théry, *Formal Proofs of Tarjan's Strongly Connected 2 Components Algorithm in Why3, Coq and Isabelle*, 2019
- [2] V. Bloemen, A. Laarman, J. van de Pol, *Multi-Core On-The-Fly SCC Decomposition*, 2016
- [3] V. Bloemen, *Strong Connectivity and Shortest Paths for Checking Models*, 2019