

Machine-checked correctness and complexity of a Union-Find implementation

Arthur Charguéraud François Pottier



December 16, 2015

Message

Let's begin with a demo...

Proving correctness and termination is not enough!

Verification methodology

We extend the **CFML** logic and tool with **time credits**.

This allows reasoning about the correctness and (amortized) complexity of realistic (imperative, higher-order) OCaml programs.

Separation Logic

Heap predicates:

$$H : \text{Heap} \rightarrow \text{Prop}$$

Usually, Heap is $\text{loc} \mapsto \text{value}$. The basic predicates are:

$$[] \quad \equiv \quad \lambda h. h = \emptyset$$

$$[P] \quad \equiv \quad \lambda h. h = \emptyset \wedge P$$

$$H_1 \star H_2 \quad \equiv \quad \lambda h. \exists h_1 h_2. h_1 \perp h_2 \wedge h = h_1 \uplus h_2 \wedge H_1 h_1 \wedge H_2 h_2$$

$$\exists x. H \quad \equiv \quad \lambda h. \exists x. H h$$

$$l \hookrightarrow v \quad \equiv \quad \lambda h. h = (l \mapsto v)$$

Separation Logic with time credits

We wish to introduce a new heap predicate:

$$\$n : \text{Heap} \rightarrow \text{Prop} \quad \text{where } n \in \mathbb{N}$$

Intended properties:

$$\$(n + n') = \$n \star \$n' \quad \text{and} \quad \$0 = []$$

Intended use:

A time credit is a permission to perform **“one step”** of computation.

Connecting computation and time credits

Idea:

- ▶ Make sure that **every function call consumes one time credit**.
- ▶ Provide **no way of creating** a time credit.

Thus,

$$(\text{total number of function calls}) \leq (\text{initial number of credits})$$

Ensuring that every call consumes one credit

The CFML tool inserts a call to `pay()` at the beginning of every function.

```
let rec find x =  
  pay();  
  match !x with  
  | Root _ -> x  
  | Link y -> let z = find y in x := Link z; z
```

The function `pay` is fictitious. It is axiomatized:

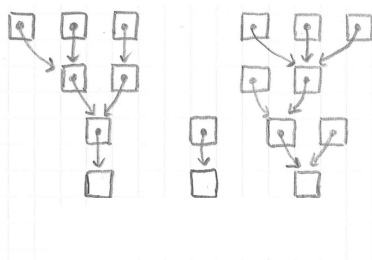
$$\text{App } \text{pay } () (\$1) (\lambda_ . [])$$

This says that **`pay()` consumes one credit.**

Contributions

- ▶ The first **machine-checked complexity analysis** of Union-Find.
- ▶ Not just at an abstract level, but **based on the OCaml code**.
- ▶ **Modular**. We establish a specification for clients to rely on.

The Union-Find data structure: OCaml interface



```
type elem
val make : unit -> elem
val find : elem -> elem
val union : elem -> elem -> elem
```

The Union-Find data structure: OCaml implementation

Pointer-based, with path compression and union by rank:

```
type rank = int

type elem = content ref

and content =
  | Link of elem
  | Root of rank

let make () = ref (Root 0)

let rec find x =
  match !x with
  | Root _ -> x
  | Link y ->
      let z = find y in
      x := Link z;
      z

let link x y =
  if x == y then x else
  match !x, !y with
  | Root rx, Root ry ->
      if rx < ry then begin
        x := Link y;
        y
      end else if rx > ry then begin
        y := Link x;
        x
      end else begin
        y := Link x;
        x := Root (rx+1);
        x
      end
  | _, _ -> assert false

let union x y = link (find x) (find y)
```

Complexity analysis

Tarjan, 1975: the **amortized** cost of union and find is $O(\alpha(N))$.

- ▶ where N is a fixed (pre-agreed) bound on the number of elements.

Streamlined proof in *Introduction to Algorithms*, 3rd ed. (1999).

$$A_0(x) = x + 1$$

$$\begin{aligned} A_{k+1}(x) &= A_k^{(x+1)}(x) \\ &= A_k(A_k(\dots A_k(x)\dots)) \quad (x + 1 \text{ times}) \end{aligned}$$

$$\alpha(n) = \min\{k \mid A_k(1) \geq n\}$$

Quasi-constant cost: for all practical purposes, $\alpha(n) \leq 5$.

Specification of `find`

Theorem `find_spec` : $\forall N D R x, x \in D \rightarrow$
App `find` `x`
 $(\text{UF } N D R \star \$(\text{alpha } N + 2))$
 $(\text{fun } r \Rightarrow \text{UF } N D R \star \backslash[r = R x]).$

The **abstract predicate** $\text{UF } N D R$ is the invariant.
It asserts that the data structure is well-formed and that we own it.

- ▶ D is the set of all elements, i.e., the domain.
- ▶ N is a bound on the cardinality of the domain.
- ▶ R maps each element of D to its representative.

Specification of union

Theorem `union_spec` : $\forall N D R x y, x \in D \rightarrow y \in D \rightarrow$
App `union` `x y`
 $(\text{UF } N D R \star \$ (3 * (\text{alpha } N) + 6))$
 $(\text{fun } z \Rightarrow$
 $\text{UF } N D (\text{fun } w \Rightarrow \text{If } R w = R x \vee R w = R y \text{ then } z \text{ else } R w)$
 $\star [z = R x \vee z = R y])$.

The amortized cost of `union` is $3\alpha(N) + 6$.

Definition of Φ , on paper

$$\begin{aligned} p(x) &= \text{parent of } x && \text{if } x \text{ is not a root} \\ k(x) &= \max\{k \mid K(p(x)) \geq A_k(K(x))\} && \text{(the **level** of } x\text{)} \\ i(x) &= \max\{i \mid K(p(x)) \geq A_{k(x)}^{(i)}(K(x))\} && \text{(the **index** of } x\text{)} \\ \phi(x) &= \alpha(N) \cdot K(x) && \text{if } x \text{ is a root or has rank 0} \\ \phi(x) &= (\alpha(N) - k(x)) \cdot K(x) - i(x) && \text{otherwise} \\ \Phi &= \sum_{x \in D} \phi(x) \end{aligned}$$

For some intuition, see [Seidel and Sharir \(2005\)](#).

Definition of Φ , in Coq

Definition p F x :=

epsilon (fun y \Rightarrow F x y).

Definition k F K x :=

Max (fun k \Rightarrow K (p F x) \geq A k (K x)).

Definition i F K x :=

Max (fun i \Rightarrow K (p F x) \geq iter i (A (k F K x)) (K x)).

Definition phi F K N x :=

If (is_root F x) \vee (K x = 0)

then (alpha N) * (K x)

else (alpha N - k F K x) * (K x) - (i F K x).

Definition Phi D F K N :=

Sum D (phi F K N).

Machine-checked amortized complexity analysis

Proving that the invariant is preserved naturally leads to this goal:

$$\Phi + \text{advertised cost} \geq \Phi' + \text{actual cost}$$

For instance, in the case of `find`, we must prove:

$$\Phi(F) + (\alpha N + 2) \geq \Phi(F') + (d + 1)$$

where:

- ▶ F is the graph before the execution of `find x`,
- ▶ F' is the graph after the execution of `find x`,
- ▶ d is the length of the path in F from x to its root.

Summary

- ▶ A machine-checked proof of **correctness and complexity**.
- ▶ Down to the level of the **OCaml code**.
- ▶ **3000 loc** of high-level mathematical analysis.
- ▶ **400 loc** of specification and low-level verification.
- ▶ Future work: write $O(\alpha(n))$ instead of $3\alpha(n) + 6$.

<http://gallium.inria.fr/~fpottier/dev/uf/>