

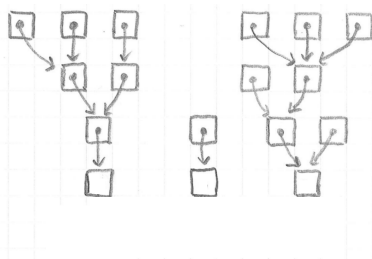
Machine-checked correctness and complexity of a Union-Find implementation

Arthur Charguéraud François Pottier



September 8, 2015

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, $\alpha(n) \leq 5$.

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.

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.

Space of the related work:

- Verification that ignores complexity.
- Verification that includes complexity:
 - Proof only at an abstract mathematical level.
 - Proof that goes down to the level of the source code:
 - with emphasis on automation (e.g., the RAML project);
 - with emphasis on expressiveness (Atkey; this work).

Specification

Separation Logic with time credits

Union-Find: invariants

Conclusion

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`
`(UF N D R * $(3*(alpha N)+6))`
`(fun z =>`
`UF N D (fun w => If R w = R x v R w = R y then z else R w)`
`* [z = R x v z = R y]).`

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

- ▶ Reasoning with O 's is ongoing work.
- ▶ Asserting that the worst-case cost is $O(\log N)$ would require non-storable time credits.

Specification of make

Theorem `make_spec` : $\forall N D R, \text{card } D < N \rightarrow$

`App make tt`

`(UF N D R * $1)`

`(fun x \Rightarrow UF N (D \cup {x}) R * \backslash [x \notin D] * \backslash [R x = x]).`

The cost of `make` is $O(1)$.

At most N elements can be created.

Specification of the ghost operations

Theorem `UF_create` : $\forall N,$
 $\llbracket \rrbracket \triangleright (\text{UF } N \ \emptyset \ \text{id}).$

Theorem `UF_properties` : $\forall N \ D \ R, \text{UF } N \ D \ R \triangleright \text{UF } N \ D \ R \ \star$
 $[(\text{card } D \leq N) \wedge$
 $\forall x, (R \ (R \ x) = R \ x) \wedge$
 $(x \in D \rightarrow R \ x \in D) \wedge$
 $(x \notin D \rightarrow R \ x = x)].$

`UF_create` initializes an empty Union-Find data structure.
It can be thought of as a ghost operation. N is fixed at this moment.

`UF_properties` reveals a few properties of D , N and R .

Specification

Separation Logic with time credits

Union-Find: invariants

Conclusion

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.

Model of time credits

We change Heap to $(\text{loc} \mapsto \text{value}) \times \mathbb{N}$.

A heap is a (partial) memory paired with a (partial) number of credits.

The predicate $\$n$ means that we own (exactly) n credits:

$$\$n \equiv \lambda(m, c). m = \emptyset \wedge c = n$$

Separating conjunction distributes the credits among the two sides:

$$(m_1, c_1) \uplus (m_2, c_2) \equiv (m_1 \uplus m_2, c_1 + c_2)$$

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,

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

This, we prove (on paper).

Connecting computation and time credits

This is a formal statement of the previous claim.

Theorem (Soundness of characteristic formulae with time credits)

$$\forall mc. \left\{ \begin{array}{l} \llbracket t \rrbracket H Q \\ H(m, c) \end{array} \right. \Rightarrow \exists nvm'c'm''. \left\{ \begin{array}{l} t/m \Downarrow^n v/m' \oplus m'' \\ n \leq c - c' \\ Qv(m', c') \end{array} \right.$$

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.**

Connecting computation and time credits

Hypotheses:

- ▶ No loops in the source code. (Translate them to recursive functions.)
- ▶ The compiler turns a function into **machine code with no loop**.
- ▶ A machine instruction executes in constant time.

Thus,

$$\begin{aligned}(\text{total \#instructions executed}) &= O(\text{total \#function calls}) \\(\text{total execution time}) &= O(\text{total \#function calls}) \\(\text{total execution time}) &= O(\text{initial \#credits})\end{aligned}$$

This, we do not prove.

(It would require modeling the compiler and the machine.)

Expressive power

An assertion $\$n$ can appear in a precondition, a postcondition, a data structure invariant, etc.

That is, time credits can be **passed** from caller to callee (and back), and can be **stored** for later use.

This allows **amortized time complexity** analysis.

Specification

Separation Logic with time credits

Union-Find: invariants

Conclusion

Invariant #1: math

Definition $\text{Inv } N D F K R :=$
confined $D F \wedge$
functional $F \wedge$
 $(\forall x, \text{path } F x (R x) \wedge \text{is_root } F (R x)) \wedge$
 $(\text{finite } D) \wedge$
 $(\text{card } D \leq N) \wedge$
 $(\forall x, x \notin D \rightarrow K x = 0) \wedge$
 $(\forall x y, F x y \rightarrow K x < K y) \wedge$
 $(\forall r, \text{is_root } F r \rightarrow 2^{(K r)} \leq \text{card}(\text{descendants } F r)).$

The relation F is the graph (i.e., the disjoint set forest).

K maps every element to its rank.

D, N, R are as before.

Invariant #2: memory

CFML describes a **region** as `GroupRef M`, where the partial map M maps a memory location to the content of the corresponding memory cell.

Invariant #3: connecting math and memory

We must express the connection between M and our D, N, R, F, K .

Definition Mem $D F K M :=$

```
(dom M = D)
^ (forall x, x in D ->
  match M[x] with
  | Link y => F x y
  | Root k => is_root F x ^ k = K x
  end).
```

M contains less information than D, N, R, F, K . E.g.,

- ▶ N is ghost state;
- ▶ the rank $K(x)$ of a non-root node x is ghost state.

Invariant #4: potential

At every time, we store Φ time credits. (Φ is defined in a few slides.)

Φ depends on D, F, K, N , so the Coq invariant is $\$ (\Phi D F K N)$.

Invariants #1-#4 together

The **abstract predicate** that appears in the public specification:

Definition $UF\ N\ D\ R := \exists F\ K\ M,$
 $\backslash [Inv\ N\ D\ F\ K\ R] \star$
 $(GroupRef\ M) \star$
 $\backslash [Mem\ D\ F\ K\ M] \star$
 $\$(Phi\ D\ F\ K\ N).$

Definition of Φ , on paper

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

Don't ask... 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).

Non-constructive operators: epsilon, Max, If, Sum. Convenient!

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:

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

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.

Specification

Separation Logic with time credits

Union-Find: invariants

Conclusion

Summary

- ▶ A machine-checked proof of **correctness and complexity**.
- ▶ Down to the level of the **OCaml code**.
- ▶ **3Kloc** of high-level mathematical analysis.
- ▶ **0.4Kloc** of specification and low-level verification.

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

Future work

- ▶ Establish a **local bound** of $\alpha(n)$ instead of $\alpha(N)$ where N is fixed.
 - ▶ Follow [Alstrup et al. \(2014\)](#).
- ▶ Introduce **O notation** and write $O(\alpha(n))$ instead of $3\alpha(n) + 6$.
- ▶ Attach a **datum** to every root. Offer a few more operations.
- ▶ Develop a **verified OCaml library** of basic algorithms and data structures (with Filliâtre and others).

Appendix

The CFML approach

```
(** UnionFind.ml **)
```

```
let rec find x =  
  ...
```

```
(** UnionFind_ml.v **)
```

```
Axiom find : Func.
```

```
Axiom find_cf :  $\forall x$  H Q,  
  (...)  $\rightarrow$  App find x H Q.
```

```
(** UnionFind_proof.v **)
```

```
Theorem find_spec :  $\forall x \in D$ ,  
  App find x (...) (...).
```

```
Proof.
```

```
  intros. apply find_cf.
```

```
  ...
```

```
Qed.
```

Characteristic formulae

The characteristic formula of a term t , written $\llbracket t \rrbracket$, is a predicate such that:

$$\forall H Q. \llbracket t \rrbracket H Q \Rightarrow \{H\} t \{Q\}$$

In any state satisfying H , t terminates on v , in a state satisfying $Q v$.

Example definition:

$$\llbracket t_1 ; t_2 \rrbracket \equiv \lambda H Q. \exists H'. \llbracket t_1 \rrbracket H (\lambda _ . H') \wedge \llbracket t_2 \rrbracket H' Q$$

Characteristic formulae: sound and complete, follow the structure of the code (compositional and linear-sized), and support the frame rule.

Characteristic formula generation

$$\llbracket v \rrbracket = \lambda H Q. H \triangleright Q v$$

$$\llbracket t_1 ; t_2 \rrbracket = \lambda H Q. \exists Q'. \llbracket t_1 \rrbracket H Q' \wedge \llbracket t_2 \rrbracket (Q' tt) Q$$

$$\llbracket \text{let } x = t_1 \text{ in } t_2 \rrbracket = \lambda H Q. \exists Q'. \llbracket t_1 \rrbracket H Q' \wedge \forall x. \llbracket t_2 \rrbracket (Q' x) Q$$

$$\llbracket f v \rrbracket = \lambda H Q. \text{App } f v H Q$$

$$\llbracket \text{let } f = \lambda x. t_1 \text{ in } t_2 \rrbracket = \lambda H Q. \forall f. P \Rightarrow \llbracket t_2 \rrbracket H Q$$

$$\text{where } P = (\forall x H' Q'. \llbracket t_1 \rrbracket H' Q' \Rightarrow \text{App } f x H' Q')$$

App has type:

$$\forall A B. \text{Func} \rightarrow A \rightarrow (\text{Heap} \rightarrow \text{Prop}) \rightarrow (B \rightarrow \text{Heap} \rightarrow \text{Hprop}) \rightarrow \text{Prop}.$$

Other amortized analyses using CFML with credits

Resizable arrays

- ▶ push and pop at back in $O(1)$.

Random-access lists

- ▶ push and pop at head in $O(1)$, get and set in $O(\log n)$.

Bootstrapped chunked sequence

- ▶ push and pop at the two ends in $O(1)$, split and join in $O(B \log_B n)$.