

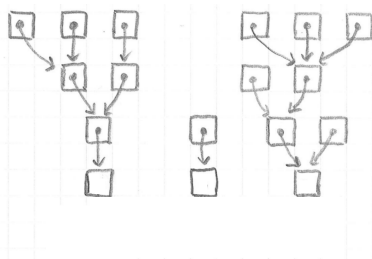
Machine-Checked Verification of the Correctness and Amortized Complexity of an Efficient Union-Find Implementation

Arthur Charguéraud
joint work with François Pottier

Inria

2015/08/25

Union-Find data structure



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

Implementation

Pointer-based Union-Find, 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)
```

Union-Find analysis

Tarjan (1975): the amortized cost of find is $O(\alpha(n))$.

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

$$A_0(x) \equiv x + 1$$

$$A_{k+1}(x) \equiv A_k^{(x+1)}(x) = A_k(A_k(\dots A_k(x)\dots)) \quad (x + 1 \text{ times})$$

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

→ In this work: the first mechanized complexity analysis of Union-Find.

Following proof from 1999, published in *Introduction to Algorithms*, 3rd ed.

Verification tool

We extend the CFML tool with time credits, to allow for the formalization of amortized complexity analyses for arbitrarily-complex OCaml programs.

Design space:

- ▶ Verification ignoring the complexity.
- ▶ Verification including the complexity:
 - ▶ Proof only at the level of the mathematical abstractions.
 - ▶ Proof also connecting to the source code:
 - ▶ with emphasis on automation (e.g., RAML project);
 - ▶ with emphasis on expressiveness (Atkey and this work).

Contents of the talk

1. Statement of specifications.
2. Separation Logic with time credits.
3. Characteristic formulae with time credits.
4. Invariant and potential for Union-Find.
5. Verification proofs.

Specification of find

Theorem find_spec : $\forall N D R x, x \in D \rightarrow$

App find x

(UF N D R \star $\$(\text{alpha } N + 2)$)

(fun r \Rightarrow UF N D R \star $\backslash[r = R x]$).

- ▶ D is the set of all elements, i.e. the domain.
- ▶ N is a bound on the cardinality of the domain.
- ▶ R maps elements to their corresponding roots.
- ▶ “UF $N D R$ ” denotes the invariant on the state.

where App has type:

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

Separation Logic

Heap predicates:

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

Core definitions:

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

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

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

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

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

→ Formalization in Coq following that of Ynot (Chlipala et al, 2009).

Principle of time credits

Time credits:

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

Properties:

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

Principle:

Ensure that every beta-reduction forces the spending of \$1.

Time credits are received in preconditions, are spent on function calls. They may be stored in the heap for later retrieval and consumption.

Requires a complexity-preserving compiler:

$$\text{nb machine instructions} = O(\text{nb beta-reductions})$$

Model of time credits

Without credits:

$$\text{Heap} \equiv (\text{loc} \mapsto \text{value})$$

With credits:

$$\text{Heap} \equiv (\text{loc} \mapsto \text{value}) \times \mathbb{N}$$

Definition of credits:

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

Earlier definitions are lifted to pairs, e.g.:

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

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.

Time credits in characteristic formulae

Goal: ensure that every beta-reduction forces the spending of \$1.

Solution: CFML instruments the OCaml code by inserting a call to “pay” at the head of every function or loop body.

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

Axiomatic specification of pay:

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

Soundness

Theorem (Soundness of characteristic formulae with time credits)

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

Union-Find invariant

Definition $\text{is_root } F \ x := \forall y, \neg F \ x \ y.$

Definition $\text{Inv } N \ D \ F \ K \ R :=$

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

- ▶ F describes the edges of the underlying graph.
- ▶ K gives the rank of every element.
- ▶ D describes the domain of the elements.
- ▶ R maps elements to their corresponding roots.
- ▶ N is a bound on the cardinality of the domain.

Representation predicate

- ▶ M maps elements to the contents of the corresponding memory cell.

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

Definition $Mem\ D\ F\ K\ M :=$
 $(dom\ M = D)$
 $\wedge (\forall x, x \in D \rightarrow$
 $match\ M[x]\ with$
 $| Link\ y \Rightarrow F\ x\ y$
 $| Root\ k \Rightarrow is_root\ F\ x \wedge k = K\ x$
 $end).$

Definition of the potential on paper

$p(x) \equiv$ parent of x (when x is not a root)

$K(x) \equiv$ rank of x

$k(x) \equiv \max\{k \mid K(p(x)) \geq A_k(K(x))\}$

$i(x) \equiv \max\{i \mid K(p(x)) \geq A_{k(x)}^{(i)}(K(x))\}$

$\phi(x) \equiv \alpha(N) \cdot K(x)$ if x is a root or has rank 0

$\phi(x) \equiv (\alpha(N) - k(x)) \cdot K(x) - i(x)$ otherwise

$$\Phi \equiv \sum_{x \in D} \phi(x)$$

Definition of the potential in Coq

Definition $p\ F\ x := \text{epsilon } (\text{fun } y \Rightarrow F\ x\ y)$.

Definition $k\ F\ K\ x := \text{Max } (\text{fun } k \Rightarrow K\ (p\ F\ x) \geq A\ k\ (K\ x))$.

Definition $i\ F\ K\ x := \text{Max } (\text{fun } i \Rightarrow K\ (p\ F\ x) \geq \text{iter } i\ (A\ (k\ F\ K\ x))\ (K\ x))$.

Definition $\text{phi } F\ K\ N\ x :=$

 If $(\text{is_root } F\ x) \vee (K\ x = 0)$

 then $(\text{alpha } N) * (K\ x)$

 else $(\text{alpha } N - k\ F\ K\ x) * (K\ x) - (i\ F\ K\ x)$.

Definition $\text{Phi } D\ F\ K\ N := \text{Sum } D\ (\text{phi } F\ K\ N)$.

Amortized analysis

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

In the case of `find`, we prove:

$$\Phi_D(F, K, N) + (\alpha N + 2) \geq \Phi_D(F', K, N) + (d + 1)$$

where:

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

Verification script

Theorem find_spec : $\forall N D R, x \in D \rightarrow$

App find x (UF N D R * \$(alpha N + 2)) (fun r \Rightarrow UF N D R * \setminus [R x = r]).

Proof.

asserts S': ($\forall d D R F K F' M,$

Inv N D F K R \rightarrow Mem D F K M $\rightarrow x \in D \rightarrow$ ipc F x d F' \rightarrow

App find x (GroupRef M * \$(d+1))

(fun r' $\Rightarrow \exists M',$ GroupRef M' * \setminus [Mem D F' K M' \wedge r' = R x])).

{ xinduction_heap Wf_nat.lt_wf. apply find_cf.

intros x d IH. hide IH. introv HI HM Dx HC. credits_split. xpay.

lets HMD: (Mem_dom HM). unfold elem in *.

xapps*. forwards* HV: (Mem_val HM) x. unfold elem in *.

xmatch; rename H0 into HK; rewrite \leftarrow HK in HV.

(* case root *) (* ... 3 lines not shown ...*)

(* case link *) (* ... 14 lines not shown ...*)

}

xweaken S'. clear S'. simpl. intros x. intros RS LRS KRS.

introv Dx. unfold UF. xextract as F K M HI HM.

forwards* (d&F'&HC&HP): amortized_cost_of_iterated_path_compression x N.

change (2%nat) with (1+1)%nat. forwards (H'&E): credits_nat_le_rest HP.

xchange E. chsimpl. credits_split. xgc H'. xframe (\$ Phi D F' K N).

xapply (>> KRS HI HM Dx HC). change (S d) with (1+d)%nat. chsimpl.

xok. clears RS. intros z. xextract as M' (HM'&Hz). subst z. intros. hsimpr~.

constructor; eauto using is_rdsf_bw_ipc, bw_ipc_preserves_RF_agreement.

Qed.

Future work

- State bounds using $\alpha(n)$, instead of $\alpha(N)$ with the constraint $n \leq N$.
- Improve the degree of inference and automation.
- Introduce the big- O notation, to write $O(\alpha(n))$ instead of $3\alpha(n) + 6$.

Conclusion

An integrated verification framework for proving not just the functional correctness but also the asymptotic complexity of a concrete program.

TCB: classic Coq + CFML generator + a complexity-preserving compiler.

Application to the verification of an efficient Union-Find implementation.

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

Thanks!