

CPSC 311: State (DRAFT) ("lec-state")

Joshua Dunfield

University of British Columbia

November 18, 2016

1 State

All of our Fun dialects have had only *immutable* bindings: During evaluation, once an identifier is bound to an expression, its meaning cannot change—it will have the same meaning as long as it is in scope (and not shadowed by another binding).

1.1 Classifying languages

Many languages have *mutable* state in some form:

- *By default, and idiomatic*: Fortran, Algol-60, Lisp, C, C++, Java, Smalltalk, . . .
- *By default, but less idiomatic*: Racket
- *Not by default*: Standard ML, OCaml
- *By simulation*: Haskell

The line between “functional” and “imperative” is fuzzy, but I think most people would draw it somewhere around Racket. The line between “purely functional” and “impurely functional”—*purity* meaning a “lack of side effects (such as state)”—is usually drawn between ML and Haskell. That line is also subject to debate, however.

Starting from the top of the list, in languages like Java, most features are mutable by default (unless `const` is given).

Racket occupies a strange position in this space: fundamental binding operations like **define** and **let** are mutable, but “good Racket style” discourages you from exploiting this. A few language features in Racket, including lists, are genuinely immutable by default (Racket also has mutable lists, but not by default; as we saw in our discussion of classifying languages, different languages often provide the same behaviours and differ only in which behaviour is the default).

The ML languages are fairly consistent in being immutable by default. A value bound by a `let` in SML or OCaml cannot be mutated. Both languages do have features similar to Racket’s boxes, but these must be used explicitly; the default is immutability. An exception that puts OCaml slightly nearer the top of the page: strings are mutable, as you can see from the following interaction with OCaml.

```
# let s = "abcd" ;;
val s : string = "abcd"
# String.set s 2 'r' ;;
- : unit = ()
# s ;;
- : string = "abrd"
```

(2016 update: I believe this has been changed in the latest version of OCaml.)

Haskell is usually considered “pure” or “purely functional”, though there is debate about this too, partly because some people argue that nontermination is a side effect. In practice, Haskell has ample support (features like the appallingly named “monads”) for an imperative style of programming. (As an aside, the techniques used to *implement* Haskell depend on mutable state!)

1.2 Defining state

The particular form of state we’ll add to Fun is called *references* (using the ML terminology). A *reference* or *ref* is essentially a pointer to a *cell* (or “ref cell”) whose contents can be mutated.

Modelling refs in a dynamic semantics requires significant changes. Now that we have an environment-based semantics, the environment env might seem to be a logical place to store the current contents of ref cells. That turns out to be a bad idea—we want the state of ref cells to survive a lexical scope, but not the state of binders (see the question below)—so instead we’ll create a new animal, a *store* S . Like an environment, a store is basically a list:

Stores $S ::= \emptyset$ empty store
 $| \ell \triangleright v, S$ location ℓ points to cell with contents v , followed by store S

Let’s extend the concrete and abstract syntax.

Expressions $\langle E \rangle ::= \dots$
 $| \{\text{Ref } \langle E \rangle\}$
 $| \{\text{Deref } \langle E \rangle\}$
 $| \{\text{Setref } \langle E \rangle \langle E \rangle\}$

The intended semantics is:

- $\{\text{Ref } \langle E1 \rangle\}$ evaluates $\langle E1 \rangle$ and returns the location of a new cell (think of a location as a pointer);
- $\{\text{Deref } \langle E1 \rangle\}$ evaluates $\langle E1 \rangle$, which must evaluate to a location, and returns the contents of the cell at that location;
- $\{\text{Setref } \langle E1 \rangle \langle E2 \rangle\}$ evaluates $\langle E1 \rangle$, which must evaluate to a location ℓ , then evaluates $\langle E2 \rangle$ and puts that resulting value into the cell at location ℓ .

The abstract syntax is extended correspondingly:

```
(define-type E
  ...
  [Ref (initial-contents E?)]
  [Deref (loc-expr E?)]
  [Setref (loc-expr E?) (new-contents E?)])
```

§1 State

We aren't quite done, though: a Ref is how we *create* a new cell, not a *pointer* to a new cell. So we need one more variant:

```
(define-type E
  ...
  [Ref (initial-contents E?)]
  [Deref (loc-expr E?)]
  [Setref (loc-expr E?) (new-contents E?)]
  [Location (locsymb symbol?)])
```

Racket has a built-in function called `gensym` that we'll use when we need a new location, to get a "fresh" symbol.

The point of a store is that its contents are mutable: evaluating an expression may change the store. So we need both an *input* store and an *output* store.

$\boxed{env; S \vdash e \Downarrow v; S'}$ Starting in environment env and store S , evaluating e produces value v and updated store S'

$$\frac{env; S \vdash e \Downarrow v; S1}{env; S \vdash (Ref\ e) \Downarrow (Location\ \ell); \ell \triangleright v, S1} \quad ??SEnv-ref$$

What is ℓ ? It really doesn't matter, as long as it isn't already in $S1$. The judgment " ℓ fresh for $S1$ " means that ℓ is not already mapped by $S1$.

$$\frac{\ell\ \text{fresh for } S1 \quad env; S \vdash e \Downarrow v; S1}{env; S \vdash (Ref\ e) \Downarrow (Location\ \ell); \ell \triangleright v, S1} \quad SEnv-ref$$

When a new feature (like Ref) leads us to change the judgment form, we need to check two things:

- We can write rules for the new features.
- We can update the rules for old features.

We seem to have a rule for the new feature Ref, so we should try to update an old rule. We'll do the rule for Pair. We first need to update Eval-pair to use environments. Since pairs don't bind identifiers, this is straightforward:

$$\frac{env \vdash e1 \Downarrow v1 \quad env \vdash e2 \Downarrow v2}{env \vdash (Pair\ e1\ e2) \Downarrow (Pair\ v1\ v2)} \quad Env-pair\ (stateless\ version)$$

$$\frac{env; S \vdash e1 \Downarrow v1; S1 \quad env; S1 \vdash e2 \Downarrow v2; S2}{env; S \vdash (Pair\ e1\ e2) \Downarrow (Pair\ v1\ v2); S2} \quad SEnv-pair$$

This version of Env-pair says that, to evaluate a pair, we first evaluate $e1$ under the given store S , producing a (possibly) changed store $S1$; then, we evaluate $e2$ under $S1$, producing another store $S2$, which is the store produced by the entire evaluation of $(Pair\ e1\ e2)$.

Following this pattern of passing stores along from premise to premise means that when we draw a derivation tree, and draw a line (really a curve) from the conclusion's starting store (to the

left of the turnstile \vdash) to the conclusion's result (to the right of the semicolon), the line looks like a thread. The store is “threaded through” the derivation tree.

Unlike previous evaluation rules for Pair, SEnv-pair specifies that an interpreter must evaluate the two expressions e_1 and e_2 in a particular order. The expression e_2 *cannot* be evaluated before e_1 , because we need to evaluate e_1 to know what S_1 is.

(Adding the full, tedious set of error-handling rules to the old big-step semantics, or to the environment-based semantics without state, would also enforce this order. But now, it is enforced within the non-error rule. Small-step semantics also enforces this, using evaluation contexts \mathcal{C} .)

Before we try to update all the old rules for this different evaluation judgment, we should figure out how to evaluate the other new features:

$$\frac{\ell \text{ fresh for } S_1 \quad \text{env}; S \vdash e \Downarrow v; S_1}{\text{env}; S \vdash (\text{Ref } e) \Downarrow (\text{Location } \ell); \ell \triangleright v, S_1} \text{SEnv-ref}$$

$$\frac{\text{env}; S \vdash e \Downarrow (\text{Location } \ell); S_2 \quad \text{lookup-loc}(S_2, \ell) = v}{\text{env}; S \vdash (\text{Deref } e) \Downarrow v; S_2} \text{SEnv-deref}$$

$$\frac{\text{env}; S \vdash e_1 \Downarrow (\text{Location } \ell); S_1 \quad \text{env}; S_1 \vdash e_2 \Downarrow v_2; S_2 \quad \text{update-loc}(S_2, \ell, v_2) = S_3}{\text{env}; S \vdash (\text{Setref } e_1 \ e_2) \Downarrow v_2; S_3} \text{SEnv-setref}$$

The idea of *update-loc* is that $\text{update-loc}(S_2, \ell, v_2) = S_3$ where S_3 is the same as S_2 , but with $\ell \triangleright \dots$ replaced by $\ell \triangleright v_2$. For example, if

$$S_2 = \ell_1 \triangleright (\text{Num } 1), \ell_2 \triangleright (\text{Num } 0), \emptyset$$

then

$$\text{update-loc}(S_2, \ell_2, (\text{Num } 2)) = \ell_1 \triangleright (\text{Num } 1), \ell_2 \triangleright (\text{Num } 2), \emptyset$$

Question: Why do we need a separate store? What goes wrong if we use the environment to store ref cells?

Because then we would end up with a similar problem as not doing a freeness check during substitution: we could access an identifier that should be out of scope.

$$e_{\text{Pair}} = (\text{Pair } (\text{Let } x \ (\text{Num } 1) \ (\text{Id } x)) \ (\text{Let } y \ (\text{Num } 2) \ (\text{Id } x)))$$

Here, the second instance $(\text{Id } x)$ is not in the scope of $(\text{Let } x \dots)$. But (if the environment contains ref cells), we have to thread the environment through; otherwise, the second component of the following pair $e_{\text{Pair}'}$ would be unable to see the effects of the first component. We expect $e_{\text{Pair}'}$ to evaluate to $(\text{Pair } (\text{Num } 1) \ (\text{Num } 1))$ because the first component changes the contents of r from $(\text{Num } 0)$ to $(\text{Num } 1)$:

$$e_{\text{Pair}'} = \left(\text{Let } r \ (\text{Ref } (\text{Num } 0)) \ (\text{Pair } (\text{Let } x \ (\text{Num } 1) \ (\text{Setref } r \ (\text{Num } 1))) \ (\text{Let } y \ (\text{Num } 2) \ (\text{Deref } r))) \right)$$

If we thread the environment through, then in e_{Pair} , the binding of $(\text{Id } x)$ would survive and could be used outside its scope; if we don't thread the environment through, $e_{\text{Pair}'}$ wouldn't behave as expected.

1.3 First implementation: `env-state.rkt`

We can **define-type** `Store` following the pattern of `Env`, and update `env-interp` to take *and* return a store. It's quite irritating, because Racket doesn't provide great support for returning pairs of things—I had to **define-type** `Config` to represent both an expression—the value v being returned in $\text{env}; S1 \vdash e \Downarrow v; S2$ —and the output store $S2$. But it can be done, and there are no big surprises.

Some of this could be done more easily in ML or Haskell, because those languages have more general “pattern matching” than **type-case**, so adjacent **type-cases** can be combined in one. We still have to look up locations in the store, and update a cell by constructing a new store with different contents for that one cell.

■ **Question:** Racket has boxes! Why not just use those, instead of going to all this trouble?

Good question. One answer is that we want to know that our interpreter follows the rules (is “sound with respect to the rules”). The code is annoying to read, but the correspondence to the rules is clear. If we use Racket's boxes as our locations, the code becomes much simpler, but we have to trust that Racket's semantics for boxes matches our rules.

I'm pretty sure it does match the rules, which is why I wrote another version of the interpreter that *does* use Racket boxes (`env-state-direct.rkt`—the word “direct” refers to using Racket's boxes directly).

Another answer is that, while we can (I think!) use Racket's boxes to correctly represent Fun's refs, we are depending on Racket's idea of a store being the same as ours. What if we wanted to allow “time travel” in Fun, where we could “checkpoint” an old store and “rewind” to it later? Racket—as far as I know—doesn't have that feature. So we'd need to either figure out how to checkpoint Racket boxes, or use the `env-state.rkt` representation where we don't use boxes.

This leads us to...

1.4 Second implementation: `env-state-direct.rkt`

In this interpreter, the operations on the `Store` **define-type** are simply calls to Racket's `unbox` and `set-box!`. Instead of reflecting the “threading” of the store in our interpreter, we assume that Racket's store behaves in that same way. (Again, I'm pretty sure it does.)

2 Translation dictionary

(Section added for lecture, 2016-11-14.)

	Racket	Fun	C	C++	Java
allocate	<code>(box e)</code>	<code>(Ref e)</code>	<code>malloc</code>	<code>new</code>	<code>new</code>
get	<code>(unbox e)</code>	<code>(Deref e)</code>	<code>*e</code>	<code>*e</code>	<code>e.fld</code>
set	<code>(set-box! e)</code>	<code>(Setref e1 e2)</code>	<code>*e1 = e2</code>	<code>*e1 = e2</code>	<code>e1.fld = e2</code>

The Java column is approximate: a Java object has any number of instance variables, rather than just one.