

Variables and Variable Binding

CPSC 509: Programming Language Principles

Ronald Garcia*

2 November 2015

In this set of notes, we learn about *variables* and *variable bindings* in programming languages. Along the way, we learn a bunch more about the nature of recursive function definitions.

So far the programs that we have been writing essentially boil down to basic arithmetic. We can write down values and computations on those values, but we have no forms of *abstraction* at all. For example, you may find yourself writing the same expression over and over again, and realize that it's always going to be the same. It would be nice if we only had to write that expression down once and then had a way to reuse it. This is exactly what *variables* (a.k.a. *identifiers*) are for! They let us rewrite an expression like

```
(if (zero? (- (+ 3 2) 1))
    (+ 3 2)
    (- (+ 3 2) 1))
```

in a form that makes it explicit that we're repeating ourselves. Rather than write the expression `(+ 3 2)` repeatedly, we say once that the variable `x` stands for the expression, and refer to that variable wherever we mean to repeat the same notion:

```
(let ([x (+ 3 2)])
  (if (zero? (- x 1))
      x
      (- x 1)))
```

Here we *bind* the value of `(+ 3 2)` to the variable `x` and refer to that variable wherever we want the value.

So we can use name binding to avoid unnecessary duplication in two senses. In the first sense, we can abbreviate our program and say things once and only once. In the second sense, some language implementations can take advantage of this form of expression abstraction to improve runtime performance. Depending on the particular language design (as we discuss later), the implementation may be able to¹ evaluate the expression only once, cache the result somewhere, and then look up the cached value whenever it is needed. This can substantially improve the performance of programs (but not always!).

Notice that I am distinguishing between the *semantics* of this language feature (i.e., the conceptual meaning of variable binding as abstracting expressions), and the *pragmatics* of the feature (i.e., how this conception of variable binding can facilitate a nice implementation strategy). In practice, language design choices can be driven in either direction. Sometimes the desired semantics for a language feature suggest a particular implementation strategy. Other times the current state of the art of implementation (or experience with implementing prior languages)—or intrinsic limitations imposed by computability or complexity theory—may drive decisions about the semantics of language features. Here we focus on the (often less appreciated) direction from semantics to pragmatics.

1 Naïve Substitution

Before we add name binding to the language, let's consider what it means to have variable references in the body of such a binding. Consider our last example above. Intuitively, we want to use the meaning

*© 2015 Ronald Garcia.

¹or *have* to!

of $(+ 3 2)$ everywhere that the expression refers to x . Taking 5 as the meaning of $(+ 3 2)$,² we want to substitute 5 for x in the body of the **let** expression.

Let's formalize *just this concept* of substitution in the context of our language of Boolean expressions.

$$\begin{array}{l}
 t \in \text{TERM}, \quad x \in \text{VAR} \\
 t ::= \text{true} \mid \text{false} \mid \text{if } t \text{ then } t \text{ else } t \\
 \quad \mid x
 \end{array}$$

The set VAR is some infinite set of identifiers, where the only thing we care about is that we can tell one variable apart from another. The TERM x is a *variable reference*, which indicates where an *argument* will be eventually substituted. A useful way to think about variable references is that, rather than treating variables as TERMS, to think of variables as *indices* on terms. So whenever you see the term x , we *really* see it as standing for the TREE $\text{var}[x]()$, where **var** is an ATOM but it is paired with a variable $x \in \text{VAR}$.³ This distinction may seem a bit strange now, but when we get to describing variable binding using the **let** form (a.k.a. **let-binding**), it will hopefully make a bit more sense.

Now here, it's worth taking an aside to talk about *metavariables* and *object variables*. Throughout the course, we have talked about metavariables like t and n , which stand for TERMS and numbers in the object language respectively. Now we have our first object language where the language itself has a notion of "variables." We have to be clear that these are not the same things as metavariables, and the implications are important. When talking about specific object language variables, we'll use blue sans-serif font variable names like x and y . We don't *really* care about the syntactic nuances of a particular programming language, like whether you can have a numeral as the first character of a variable name, or whether you can have Unicode characters in your variable names, etc. But we need SOME way of writing examples, hence our blue object variables. Metavariables like x_1 and x_2 refer to metalanguage variables that represent object language variables like x and y . Probably the most important thing to keep in mind, which leads to no end of mistakes, is that in a given context, an individual metavariable always refers to the same thing. So $x = x$ is going to always be true. That's the funny thing about variables, they don't really vary (at least not the way those thingies in C that we (sadly also) call variables do). On the other hand, given two different metavariables, x_1 and x_2 , we don't know a priori whether they refer to the same metavariable or not. If we want them to definitely be different, we must *explicitly* say $x_1 \neq x_2$. Sometimes texts will leave that side condition out under the assumption that it's obvious from the context.⁴ In those circumstances it's really important to make sure that you notice that this may be what's going on, and if so to make a note of it to yourself as you try to understand what you are reading. Below, we'll see this come up.

We model substitution using functions from TERMS to TERMS. The strategy that we will take, though, is to define a different substitution function for *each* instance of substitution. The idea is that once you've seen one version of this function, you can see exactly how you would define any other (so you then know that each such function exists, is unique, and satisfies the equations used to describe it). We take this approach because it lets us use the tools that we currently have available to prove that these functions exist.

For instance, let's define a function $[x \mapsto 0] : \text{TERM} \rightarrow \text{TERM}$ that substitutes 0 for every instance of x in a term. Note that $[x \mapsto 0]$ is just an evocative *name* for this function, just like *eval* or *dom*. The 0 and x don't actually *do* anything here. We could replace the name with F and it would still be the same function. For instance, we expect the following equation to hold:

$$[x \mapsto 0](\text{if false then } x \text{ else } (x + 1)) = \text{if false then } 0 \text{ else } (0 + 1).$$

We define this function recursively, first in long-hand style. Here is the principle of definition by recursion for the language augmented with variables:

Proposition 1 (Principle of Definition by Recursion on terms $t \in \text{TERM}$). *Let S be a set and $s_t, s_f \in S$ be two elements,*

$$H_{if} : S \times S \times S \rightarrow S$$

²Stay tuned! We revisit this assumption below when we discuss *by-name* evaluation.

³In this context, When we refer to TREE, we are *really* referring to TREE[ATOM \cup ($\{\text{var}\} \times \text{VAR}$)], representing the VAR-indexed **var** as a pair.

⁴I've been guilty of doing this, I'll admit it! Sorry...

be a function on S , and

$$H_{var} : \text{VAR} \rightarrow S.$$

Then there exists a unique function

$$F : \text{TERM} \rightarrow S$$

such that

1. $F(\text{true}) = s_t$;
2. $F(\text{false}) = s_f$;
3. $F(\text{if } t_1 \text{ then } t_2 \text{ else } t_3) = H_{if}(F(t_1), F(t_2), F(t_3))$;
4. $F(x) = H_{var}(x)$.

The new components of the principle are marked in grey above. The H_{var} function on variables VAR is responsible for handling *all* possible variable references. Essentially this function is analogous to the s_t and s_f constants, which handle the two individual constants. Recall that earlier I mentioned that we treat variable references as VAR-indexed trees $\text{var}(x)$. Here is where we see how this point of view is played out. Mind you this approach is just a way of thinking about the idea of variables all being the same kind of term. We could do the same thing with numbers n if the language had them.

So now let's define $[x \mapsto 0]$. We define the function by applying the principle of definition by recursion to the following components:

1. $S = \text{TERM}$;
2. $s_t = \text{true}$;
3. $s_f = \text{false}$;
4. $H_{if}(t_1, t_2, t_3) = \text{if } t_1 \text{ then } t_2 \text{ else } t_3$;
5. $H_{var}(x) = 0$
 $H_{var}(x_0) = x_0$ if $x_0 \neq x$.

Be careful not to be confused by the fact that we chose $S = \text{TERM}$, which is necessary to define a function from TERMS to TERMS.

Substituting these components into the principle itself (and breaking the cases into separate equations as is standard practice), we get the following shorthand definition:⁵

$$\begin{aligned} [x \mapsto 0] : \text{TERM} &\rightarrow \text{TERM} \\ [x \mapsto 0]x &= 0 \\ [x \mapsto 0]x_0 &= x_0 \quad \text{if } x_0 \neq x \\ [x \mapsto 0]\text{true} &= \text{true} \\ [x \mapsto 0]\text{false} &= \text{false} \\ [x \mapsto 0](\text{if } t_1 \text{ then } t_2 \text{ else } t_3) &= \text{if } ([x \mapsto 0]t_1) \text{ then } ([x \mapsto 0]t_2) \text{ else } ([x \mapsto 0]t_3). \end{aligned}$$

If we extend our language to include arithmetic expressions, and follow the same recipe, then we can show that the example above exactly fits this model.

Now to generalize: one thing that we can immediately see is that given any TERM t and any VAR x , we can define a function $[x \mapsto t] : \text{TERM} \rightarrow \text{TERM}$ by choosing a different $H_{var}(x)$ function for the principle of recursion. This fact justifies a corresponding generic function definition.

$$\begin{aligned} [\cdot \mapsto \cdot] : \text{VAR} \times \text{TERM} &\rightarrow \text{TERM} \rightarrow \text{TERM} \\ [x \mapsto t]t_0 &= [x \mapsto t]t_0 \end{aligned}$$

⁵It is a common notational style to not wrap the substitution argument in parentheses.

For this definition, we're using slightly different notation (no underlines!) just to make it clear that we are defining a separate general substitution function using the hard-wired individual functions. We call this function *naïve* substitution as a hint that we are working toward a fully satisfactory notion of *capture-avoiding* substitution, which comes later. Now our approach to defining naïve substitution here may seem contrived: defining individual substitution functions first and then defining the big mother of functions later. Well to be honest, *it is!* It's possible to define naïve substitution outright as a group of recursive equations without a side step through individual substitution functions, and as you'll see the definition looks pretty familiar. However, we have not yet learned the tools that we would need to claim the legitimacy of this equational definition without also providing a proof that the equations indeed define a function (i.e., a proof along the lines of the proof of the principle of definition by recursion). Later in the course we'll learn more tools so that we can justify such a definition without writing down a separate proof. However, we know enough to be able to define substitution this way, and externally we could prove some equational theorems about it that make it easier to use directly.

Proposition 2. *Naïve substitution satisfies the following equations.*

$$\begin{aligned}
 [x \mapsto t]x &= t \\
 [x \mapsto t]x_0 &= x_0 \quad \text{if } x_0 \neq x \\
 [x \mapsto t]\mathbf{true} &= \mathbf{true} \\
 [x \mapsto t]\mathbf{false} &= \mathbf{false} \\
 [x \mapsto t](\mathbf{if } t_1 \mathbf{ then } t_2 \mathbf{ else } t_3) &= \mathbf{if } ([x \mapsto t]t_1) \mathbf{ then } ([x \mapsto t]t_2) \mathbf{ else } ([x \mapsto t]t_3).
 \end{aligned}$$

Funny, this *proposition* look almost exactly like our *definition* of the $[x \mapsto 0]$ function. Furthermore, we can prove that there is *exactly one* function in $\text{VAR} \times \text{TERM} \rightarrow \text{TERM} \rightarrow \text{TERM}$ that satisfies these equations, so this could in fact be our definition, but we have to prove that this is true before we can assert that this is a definition. Remember: the principle of definition by recursion once-and-for-all proves this uniqueness property for a certain class of function definitions. Unfortunately, this set of equations cannot be recast in a way that fits the *schema* put forth by that definition (you should ask yourself and answer: what goes wrong?). So for now we settle for our indirect way of defining the function and separately proving that the above proposition holds. Later, though, we'll be able to easily justify these equations as our definition of substitution.

2 Let Bindings

Now that we have a mathematical model for what it means to substitute a language term for a variable reference, our next step is to add **let** bindings to our language!

First, let's introduce our notation for **let** bindings.

$$t ::= \dots \mid \mathbf{let } x = t \mathbf{ in } t$$

The idea behind this expression is that **let** $x = t_1$ **in** t_2 means that within the expression t_2 we are to take the variable x to stand for the result of the expression t_1 . Referring back to tree notation, we can think of this as $\mathbf{let}[x](t, t)$, where again the variable x is just a VAR, not a TERM. Another way to think about it is that the name of the expression constructor is (**let** $x = \cdot$ **in** \cdot): it is indexed on some variable x . This corresponds to the Racket notation (**let** ($[x_1 t_1]$) t), where the extra set of parentheses around the bindings is to allow for multiple simultaneous bindings, i.e., (**let** ($[x_1 t_1] [x_2 t_2] \dots$) t_2). As a notational convention, we will sometimes use indentation to mark off the entirety of a **let** expression's body (the part after **in**).

Let's write down the rest of the formal semantics of our new language of Boolean and Let expressions (BL). The syntax is as follows:

$$\begin{aligned}
 n &\in \mathbb{Z}, & t &\in \text{TERM}, & x &\in \text{VAR} \\
 t &::= \mathbf{true} \mid \mathbf{false} \mid \mathbf{if } t \mathbf{ then } t \mathbf{ else } t \\
 & \mid \mathbf{let } x = t \mathbf{ in } t
 \end{aligned}$$

To specify our operational semantics, we'll use reduction semantics, which keeps things succinct. You should make sure that you could rephrase these semantics using structural-operational or big-step semantics as well.

We start with the necessary additional syntactic notions: values, redexes, and evaluation contexts:

$$\begin{aligned} v &\in \text{VALUE}, & r &\in \text{REDEX}, & E &\in \text{ECTXT}, \\ v &::= b \\ E &::= \square \mid E[\text{if } \square \text{ then } t \text{ else } t] \mid E[\text{let } x = \square \text{ in } t] \\ r &::= \text{if } v \text{ then } t \text{ else } t \mid \text{let } x = v \text{ in } t \end{aligned}$$

You can see from the evaluation context $E[\text{let } x = \square \text{ in } t]$ that the semantics will evaluate the *bound* expression before it processes the *body* of the **let** expression. This is consistent with $\text{let } x = v \text{ in } t$ as a REDEX. Notice that this language introduces no new VALUES to the language, just a form of named abstraction.

The notions of reduction for the language follow.

$$\begin{aligned} \rightsquigarrow &\subseteq \text{REDEX} \times \text{TERM} \\ \text{if true then } t_2 \text{ else } t_3 &\rightsquigarrow t_2 \\ \text{if false then } t_2 \text{ else } t_3 &\rightsquigarrow t_3 \\ \text{let } x = v \text{ in } t &\rightsquigarrow [x \mapsto v]t \end{aligned}$$

The only new notion of reduction is for, yup you guessed it, **let**! Once you have a value in the binding position, the semantics substitutes it into the body of the expression. A side-warning: we're not quite done yet, because we haven't said anything about how to naively substitute into a **let** expression, i.e. $[x \mapsto v]\text{let } x_1 = t_1 \text{ in } t_2$. We'll cover that in the next section.

Forging ahead, let's define our evaluator. In the past, we allowed *any* TERM to count as a program, but here we are going to make a restriction. In particular, we will not allow a program to have any references to variables that have not been bound in the surrounding expression. So an expression like: $\text{let } x = \text{true in } y$ doesn't count as a program, because our operational semantics will step this to y , and then we're stuck: there is no notion of reduction for variable references. Instead, the variable should have been replaced by some VALUE by the time we get to its position in the program. If we don't have a value for our variable when we reference it, then how can we sensibly proceed?!⁶ To capture this, we must introduce a notion of *free variables*. The free variables of an expression are the variables that are referenced in a spot where no surrounding **let** binding provides their value. We express this idea in precise form as a function:

$$\begin{aligned} FV &: \text{TERM} \rightarrow \mathcal{P}(\text{VAR}) \\ FV(\text{true}) &= \emptyset \\ FV(\text{false}) &= \emptyset \\ FV(\text{if } t_1 \text{ then } t_2 \text{ else } t_3) &= FV(t_1) \cup FV(t_2) \cup FV(t_3) \\ FV(x) &= \{x\} \\ FV(\text{let } x = t_1 \text{ in } t_2) &= FV(t_1) \cup (FV(t_2) \setminus \{x\}) \end{aligned}$$

As you can see, a **let** binding handles all free references to the bound variable within its body. Now we can explicitly define our programs as the set of *closed* terms: terms that have no free variables:

$$\begin{aligned} \text{CLOSEDTERM} &= \{t \in \text{TERM} \mid FV(t) = \emptyset\} \\ \text{PGM} &= \text{CLOSEDTERM}, & \text{OBS} &= \text{VALUE} \\ \text{eval} &: \text{PGM} \rightarrow \text{OBS} \\ \text{eval}(t) &= b \text{ if } t \longrightarrow^* b \end{aligned}$$

Our language definition is almost complete. We just have to extend substitution to handle **let** bindings. For our immediate purposes, this will be straightforward, but we'll find that extending it to the general case raises some issues.

⁶Emphasis here is on *sensibly*. I'm looking at you C and JavaScript programming languages! The Python programming language allows an unbound variable reference as part of a program, and only fails if that reference is evaluated, producing a "NameError".

3 Naïve Substitution into Let Bindings

Since our language allows us to put `let` expressions anywhere, we need to be able to handle them in the bodies of one another. A couple of examples follow:

- (1) `let x = 7 in let y = 6 in x`
- (2) `let x = 7 in let x = 6 in x`

If we apply the reduction semantics to the first example above, then we expect to get:

$$\begin{aligned} & \text{let } x = 7 \text{ in let } y = 6 \text{ in } x \\ \longrightarrow & [x \mapsto 7] \text{let } y = 6 \text{ in } x \\ = & [x \mapsto 7] \text{let } y = 6 \text{ in } x \end{aligned}$$

Similarly for the second example we get:

$$\begin{aligned} & \text{let } x = 7 \text{ in let } x = 6 \text{ in } x \\ \longrightarrow & [x \mapsto 7] \text{let } x = 6 \text{ in } x \\ = & [x \mapsto 7] \text{let } x = 6 \text{ in } x \end{aligned}$$

So what should $[x \mapsto 7] \text{let } y = 6 \text{ in } x$ be? Well, y doesn't really have anything to do with x , so the obvious thing to do is substitute for x in the body of the `let`. In the case of $[x \mapsto 7] \text{let } x = 6 \text{ in } x$ though, it sure looks like x should become `6` once the inner `let` expression is evaluated, so the proper behaviour seems to be to let it alone.⁷

More generally, we expect naïve substitution to satisfy the following equations:

$$\begin{aligned} [x \mapsto 7] \text{let } x = t_1 \text{ in } t_2 &= \text{let } x = [x \mapsto 7]t_1 \text{ in } t_2 \\ [x \mapsto 7] \text{let } x_0 = t_1 \text{ in } t_2 &= \text{let } x_0 = [x \mapsto 7]t_1 \text{ in } [x \mapsto 7]t_2 \quad x_0 \neq x \end{aligned}$$

In short, it should not mess with inner bindings of the same variable. In both cases we substitute into the expression that is to be bound to the variable though.

At this point, we might feel satisfied with our equations, but they introduce a bit of gum into the works. If we were to extend our Principle of Definition by Recursion following our pattern to date, we would simply add an extra function:

$$H_{let} : \text{VAR} \times S \times S \rightarrow S$$

And the unique function F would satisfy the equation:

$$F(\text{let } x_0 = t_1 \text{ in } t_2) = H_{let}(x_0, F(t_1), F(t_2)).$$

This raises a problem: there is no way to define H_{let} to satisfy our equations! By the time we know that $x_0 = x$, the body of the expression has already been substituted into. Shucks!

Never fear though: it turns out that our principle of definition by recursion can be generalized. The simple one that we used, which proceeds by *structural recursion*, discards the original terms. But that was never strictly necessary: it just happens to be a special case that fits the bill in the grand majority of cases. We introduced that version so as to not make things too complex from the start.

To handle our notion of naïve substitution, we need a stronger principle of recursion.

Proposition 3 (Principle of Definition by Primitive Recursion on $t \in \text{TERM}$). *Let S be a set and $s_t, s_f \in S$ be two elements,*

$$\begin{aligned} H_{if} &: \text{TERM} \times \text{TERM} \times \text{TERM} \times S \times S \times S \rightarrow S \\ H_{let} &: \text{VAR} \times \text{TERM} \times \text{TERM} \times S \times S \rightarrow S \\ H_{var} &: \text{VAR} \rightarrow S \end{aligned}$$

⁷Pun not intended, seriously!

be functions.

Then there exists a unique function

$$F : \text{TERM} \rightarrow S$$

such that

1. $F(\text{true}) = s_t$;
2. $F(\text{false}) = s_f$;
3. $F(x) = H_{\text{var}}(x)$.
4. $F(\text{if } t_1 \text{ then } t_2 \text{ else } t_3) = H_{\text{if}}(t_1, t_2, t_3, F(t_1), F(t_2), F(t_3))$;
5. $F(\text{let } x = t_1 \text{ in } t_2) = H_{\text{let}}(x, t_1, t_2, F(t_1), F(t_2))$.

Proof. Exercise for the reader. □

The principle of definition by primitive recursion provides enough additional structure to enable our component H functions to consider the original terms.

Armed with the principle of definition by *primitive* recursion, we can easily provide a case for `let` that lets us define the `[x ⇒ 7]` function we want:

$$\begin{aligned} H_{\text{let}}(x, t_1, t_2, t'_1, t'_2) &= \text{let } x = t'_1 \text{ in } t_2 \\ H_{\text{let}}(x_0, t_1, t_2, t'_1, t'_2) &= \text{let } x_0 = t'_1 \text{ in } t'_2 \quad x_0 \neq x. \end{aligned}$$

The main lesson here is that we should be careful with function definitions, and the full-stack semanticist should really know which recursion principle can be used to justify them. Before too long we'll learn the recursion principle to rule all recursion principles (it's called *well-founded recursion*). With that one, we'll even be able to directly define the general substitution function $[\cdot \mapsto \cdot]$. Sit tight!

4 By-Name Let Binding

At this point we have a perfectly fine language definition. But let's consider a slightly different point in the design space. When we defined our language above, we made a design decision without much fanfare: that the expression bound to a `let` is evaluated to a value before it is substituted. Most popular programming languages operate this way, but it's not the only option. Another option is to substitute *the entire expression* and evaluate it at every variable reference. This is sometimes called *call-by-name*, which is a little odd here because we have no functions to call. The "by-name" part goes back to the design of the Algol 60 language, and the idea was that the "name" was the entire expression. Our original design is called "by-value" since variables are bound only to the values of expressions. Let's consider what changes we would need to make to our semantics to support by-name let bindings. The main changes happen in the syntax for our extra forms:

$$\begin{aligned} E &::= \square \mid E[\text{if } \square \text{ then } t \text{ else } t] \\ r &::= \text{if } v \text{ then } t \text{ else } t \mid \text{let } x = t \text{ in } t \end{aligned}$$

First, we remove the evaluation context for `let` because we will no longer evaluate the bound expression. Then the redexes change in that any arbitrary `let` expression counts as a redex. Thus, our notion of reduction for `let` is updated.

$$\text{let } x = t_1 \text{ in } t_2 \rightsquigarrow [x \mapsto t_1]t_2$$

The significance of this style of application is that it may yield *worse* performance lead you to perform more work if you refer to the same variable many times, as in

$$\text{let } x = \text{BIGEXPRESSION} \text{ in if } x \text{ then } x \text{ else } x$$

Here *BIGEXPRESSION* gets recomputed 3 times. Under by-value it would only be computed once.

On the other hand, by-name can be more efficient if it avoids unnecessary work., for example:

```
let x = BIGEXPRESSION in if true then false else true
```

Under the by-name model, *BIGEXPRESSION* never gets run at all, while it gets run in by-value once.

Later we will be introduced to a computation model called “by-need” that tries to blend both by-value and by-name to get the best of both worlds (though it still costs you some things). In practice, production languages⁸ always use “by-need” to mimic by-name behaviour, but understanding by-name semantics is typically sufficient to reason about such languages, because they avoid features that would let you tell the difference.

5 Capture-Avoiding Substitution

In the semanticses above, both by-value and by-name, we take advantage of a significant property of programs: that we only substitute closed expressions into the bodies of `let` expressions. To see this, consider that we only consider closed expressions to be programs, and that furthermore, evaluation contexts do now allow us to evaluate under variable bindings, which are the only places where free variables could be found.

So evaluation can be defined using naïve substitution without any trouble, but what happens if we try to substitute a term with free variables into the body of a `let`. Well, things can go badly.

Consider the following expression:

```
let z = 3
in let y = z
  in let z = 2
    in y
```

and suppose that we can perform the reductions in any order. If we start with the outermost and work our way in, then everything goes fine, since it corresponds to our by-value evaluation which yields 3. But suppose we reduce the *middle* binding first.

$$\begin{array}{l} \text{let } z = 3 \\ \text{in let } y = z \\ \quad \text{in let } z = 2 \\ \quad \quad \text{in } y \end{array} \quad \longrightarrow \quad \begin{array}{l} \text{let } z = 3 \\ \text{in let } z = 2 \\ \quad \text{in } z \end{array}$$

Then no matter the order of the last two steps, we end up with 2 as our result. What gives? Well, the `z` that is bound to `y` is meant to refer to the binding of `z` on the outside, but by willy-nilly substituting `z` under another binding of the same variable, we *inadvertently capture z*. How can we avoid this? Well, one way is to *rename* the inner `z` binding to some other innocuous name, say `g`.

$$\begin{array}{l} \text{let } z = 3 \\ \text{in let } y = z \\ \quad \text{in let } g = 2 \\ \quad \quad \text{in } y \end{array} \quad \longrightarrow \quad \begin{array}{l} \text{let } z = 3 \\ \text{in let } g = 2 \\ \quad \text{in } z \end{array}$$

Now that we have renamed the inner `z` to `g`, we have nothing to worry about. So the lesson here is that substituting open terms into the body of a `let` must be careful to avoid inadvertent capture.

We can develop a richer notion of substitution that exactly addresses this issue. This new operation is called *capture-avoiding* substitution, and it is the standard notion of substitution.

First off, why would we want to do this rewrite out of order? The most important answer is that we would like to understand what kinds of changes we can make to a program that will preserve the result of *eval*. These are called *correctness-preserving program transformations*. It turns out that applying notions of reduction out of order are often correctness-preserving, and doing so can often improve the performance of a program. I bet that as a programmer you have essentially made these kinds of rewrites to your code

⁸The Haskell programming language is the most well-known of this sort.

before, based on the intuition that the program will still “do the same thing”. We would like to be able to formalize this kind of thing, and being able to perform substitutions under variable binders is a common example.⁹

So, we can see that naïve substitution doesn’t cut it in this case. That means that we need to develop a new kind of substitution that supports substituting terms with free variables into other terms. We pursue this now.

5.1 α -equivalence, or, Bound Variable Names Don’t Matter

One of the key observations above that helped us fix our errant instance of substitution is that the choice of bound variable names shouldn’t really matter in a program. For instance, the following two expressions:

```
let x = 7 in x
let y = 7 in y
```

Are for all intents and purposes the same. While useful variable names are helpful to humans when they want to understand programmer intent, the language semantics just doesn’t care which variable names you choose, so long as the variable references line up with the variable bindings in the same way.

The alternative is a programming language where you have to worry all the time about what local variable names are used within each function. That makes it hard to build correct program modules independently: every piece of code has to watch out for colliding with variable names from other modules...that’s almost as bad as using global variable names everywhere!

In the language of we have described so far this property that bound variable names don’t matter holds, and we can formalize it as an *equivalence relation*, using naïve substitution to help us.

Definition 1 (Alpha Equivalence). *Let \sim_a : TERM \times TERM be defined by the following rules:*

$$\begin{array}{c} \overline{x \sim_a x} \qquad \overline{b \sim_a b} \qquad \frac{t_{11} \sim_a t_{21} \quad t_{12} \sim_a t_{22} \quad t_{13} \sim_a t_{23}}{\text{if } t_{11} \text{ then } t_{12} \text{ else } t_{13} \sim_a \text{if } t_{21} \text{ then } t_{22} \text{ else } t_{23}} \\ \\ \frac{t_{11} \sim_a t_{21} \quad t_{12} \sim_a t_{22}}{\text{let } x = t_{11} \text{ in } t_{12} \sim_a \text{let } x = t_{21} \text{ in } t_{22}} \\ \\ \frac{\text{let } x_3 = t_{11} \text{ in } [x_1 \mapsto x_3]t_{12} \sim_a \text{let } x_3 = t_{21} \text{ in } [x_2 \mapsto x_3]t_{22} \quad x_1 \neq x_2, \quad x_3 \notin FV(t_{12}) \cup FV(t_{22})}{\text{let } x_1 = t_{11} \text{ in } t_{12} \sim_a \text{let } x_2 = t_{21} \text{ in } t_{22}} \end{array}$$

Alpha-equivalence, written \sim_a , is a binary relation on terms. The “alpha” bit is a painfully cute reference to “alphabet”, implying that the programs are equivalent up to the choice of names for bound variables. Almost all of the cases that define this relation look exactly like equality: If we throw out the last rule, then the definition *is* exactly equality. This makes it clear that identical terms are alpha-convertible. The last rule, though, is the only interesting one. Essentially it says that two **let** bindings with *different* bound variable names (i.e., $x_1 \neq x_2$) are alpha-convertible if *renaming their bound variables* to a common name that does not capture any previously-free variables (i.e., $x_3 \notin FV(t_{12}) \cup FV(t_{22})$), suffices to make the terms alpha-equivalent. Later you may see alternative definitions of the same notion of alpha-equivalence.

Now, when you look at the definition of \sim_a , it “relates” terms by reconciling the variable names bound by **let** expressions, but we are technically on the hook to confirm that this induces some notion of equivalence, i.e. that two terms are “equal enough for our purposes.” In particular, we should prove that \sim_a is in fact an *equivalence relation*.

Proposition 4 (\sim_a is an equivalence relation). *Given $t_1, t_2, t_3 \in \text{TERM}$, the following are true:*

1. $t_1 \sim_a t_1$;
2. If $t_1 \sim_a t_2$ then $t_2 \sim_a t_1$;

⁹A refactoring tool can break your code if it is not sufficiently careful with performing program rewrites underneath variable bindings!

3. If $t_1 \sim_a t_2$ and $t_2 \sim_a t_3$ then $t_1 \sim_a t_3$.

Equivalence relations arise *all the time* in mathematics. They are a mathematical way of characterizing groups of stuff that are “like one another” (i.e. are equivalent) so long as you ignore some uninteresting differences. *Equality* (or *identity*) = is also an equivalence relation, but it is the one that ignores nothing: two entities are only equal if they are truly identical: one and the same. In our case here, two TERMS are alpha-equivalent if they are the same except for the choice of variables bound by `let`. Our hope is that our choices for bound variable names should not “matter,” in the sense that it should not affect the meaning of a program, though it can affect whether some programmer understands what you are doing. As it turns out, we can prove that names don’t matter (at least not to evaluation).

Proposition 5. If $t_1 \sim_a t_2$ then $eval(t_1) = eval(t_2)$.

Proof. An exercise. □

Thus, α -equivalence is an example of a correctness-preserving transformation on programs, or what the kids might call “a legal refactoring”.

6 Generalizing Substitution

We just saw that evaluation, whose definition depends under the hood on naïve substitution, is invariant with respect to alpha equivalence. However, naïve substitution itself *does not* respect alpha equivalence. We can see this by considering the underlying problem with one of our earlier examples:

$$z \sim_a z$$

and

$$\begin{array}{c} \text{let } z = 2 \\ \text{in } y \end{array} \sim_a \begin{array}{c} \text{let } g = 2 \\ \text{in } y \end{array}$$

but

$$[y \mapsto z] \left(\begin{array}{c} \text{let } z = 2 \\ \text{in } y \end{array} \right) = \begin{array}{c} \text{let } z = 2 \\ \text{in } z \end{array} \not\sim_a \begin{array}{c} \text{let } g = 2 \\ \text{in } z \end{array} = [y \mapsto z] \left(\begin{array}{c} \text{let } g = 2 \\ \text{in } y \end{array} \right)$$

So as far as naïve substitution work, bound variable names matter! We didn’t run into problems earlier because our operational semantics only substitutes *closed terms*, a special case that causes no problems.

Proposition 6. If t_{11} is closed, $t_{11} \sim_a t_{21}$, and $t_{12} \sim_a t_{22}$, then $[x \mapsto t_{11}]t_{12} \sim_a [x \mapsto t_{21}]t_{22}$.

However, naïvely substituting *open terms*—terms that have free variables—can cause problems, as above.

The key to making substitution work for *arbitrary* terms, with or without free variables, rests with alpha equivalence. We’ve seen already that we can substitute open terms without incident if we carefully rename bound variables along the way. We can bake this technique directly into an enhanced version of substitution, and the resulting operation will respect alpha equivalence, as we desire.

In essence, these propositions say that if two programs differ only in their choice of bound variables, then their results also only differ in their bound variable names. Let’s formalize this with a new function definition:

$$\begin{aligned} [t/x] : \text{TERM} &\rightarrow \text{TERM} \\ [t/x]\text{true} &= \text{true} \\ [t/x]\text{false} &= \text{false} \\ [t/x](\text{if } t_1 \text{ then } t_2 \text{ else } t_3) &= \text{if } [t/x]t_1 \text{ then } [t/x]t_2 \text{ else } [t/x]t_3 \\ [t/x]x &= t \\ [t/x]x_0 &= x_0 \text{ if } x \neq x_0 \\ [t/x](\text{let } x_0 = t_1 \text{ in } t_2) &= [t/x](\text{let } x_1 = t_1 \text{ in } [x_0 \mapsto x_1]t_2) \quad x_0 \in \{x\} \cup FV(t) \\ [t/x](\text{let } x_0 = t_1 \text{ in } t_2) &= \text{let } x_0 = [t/x]t_1 \text{ in } [t/x]t_2 \quad x_0 \notin \{x\} \cup FV(t) \end{aligned}$$

where x_1 is the *least* variable such that $x_1 \notin \{x\} \cup FV(t) \cup FV(\text{let } x_0 = t_1 \text{ in } t_2)$

The key equations here are the last two. Notice how the second-to-last uses naïve substitution to rename the bound variable of the `let` expression to a new variable that does not already appear free. Mind you, we could be more conservative about renaming, only doing so where *strictly* necessary, so as to keep the original variable names wherever possible. Doing so would be very helpful for a human-facing implementation, but here we are primarily interested in the mathematics: *bound variable names don't matter!*

Given this definition, we can produce our final function $[\cdot/\cdot] : \text{TERM} \times \text{VAR} \times \text{TERM} \rightarrow \text{TERM}$.

We call this *capture-avoiding substitution* because it never accidentally captures a bound variable.

There are a few things worth mentioning about this particular definition. In the literature, a number of subtle points are typically glossed over because you can get pretty far with approximately the right definition, but it's good to know the finer details so that they don't come back to bite you.

Notice at the end the side condition about *least* variables. This function assumes that variables have some ordering x_0, x_1, x_2, \dots to make sure that the resulting function is indeed a function, that is, always maps its input term to a unique output. Yes, this is super-contrived, but to be honest there are some programming language implementations that essentially do something like this to deterministically choose new variable names. We may talk about that later in the course as we implement language features.

Confirming that the equations above, as written, take a little work. In short, we can fuse the last two equations (by properly coalescing their side-conditions), and then appeal to the principle of definition by primitive recursion to justify this function definition. Then we can define our general function ternary substitution function. We still cannot define capture-avoiding substitution in general from the start, since it's technically not defined by recursion over the structure of the last term.¹⁰

Capture-avoiding substitution satisfies the properties that we care about:

Proposition 7. *If t is closed then $[x \Rightarrow t]t_1 \sim_a [t/x]t_1$.*

Thus, capture-avoiding substitution agrees *enough* with naïve substitution for closed terms, so it can be used to define an evaluator function without breaking anything.

Proposition 8. *If $t_{11} \sim_a t_{21}$ and $t_{12} \sim_a t_{22}$ then $[t_{11}/x]t_{12} \sim_a [t_{21}/x]t_{22}$.*

Proof. By induction over $t_{11} \sim_a t_{22}$.¹¹ □

This proposition tells us that unlike naïve substitution, capture-avoiding substitution plays well with alpha-equivalence. So now renaming bound variables doesn't really matter (as long as you continue not to care about variable names...it matters a lot to programmers if you're trying to implement a debugger for instance)!

In the next section we take another approach to dealing with changing variable names.

6.1 Equivalence Classes and the Variable Convention

What makes equivalence relations special is that each one *partitions* a set into a collection of non-overlapping subsets: the sets of all items that are equivalent to one another. This is one way of interpreting a group of things as though they were conceptually one (big) thing: collect them all. In our particular case, alpha equivalence \sim_a partitions the set `TERM` into the set:

$$\text{TERM}/\sim_a = \{ T \in \mathcal{P}(\text{TERM}) \mid \exists t_1 \in \text{TERM}. t_1 \in T \wedge \forall t_2 \in \text{TERM}. t_2 \in T \iff t_1 \sim_a t_2 \}$$

of all α -equivalence classes. We'll refer to these classes by first defining a function that maps a term to its alpha-equivalence class:

$$\begin{aligned} [\cdot]_{\sim_a} &: \text{TERM} \rightarrow \text{TERM}/\sim_a \\ [t]_{\sim_a} &= \{ t_0 \in \text{TERM} \mid t \sim_a t_0 \} \end{aligned}$$

¹⁰...despite what a number of textbooks say (sigh).

¹¹I recommend working through this.

Then any equivalence class can be denoted by mapping one of the members of the class, by “bracketing” it. For example,

$$\begin{aligned} [\text{let } x_3 = 7 \text{ in } x_3]_{\sim_a} &= \{ \text{let } x_0 = 7 \text{ in } x_0, \text{let } x_1 = 7 \text{ in } x_1, \text{let } x_2 = 7 \text{ in } x_2, \dots \} \\ [\text{let } x_3 = 7 \text{ in } x_1 x_3]_{\sim_a} &= \{ \text{let } x_0 = 7 \text{ in } x_1 x_0, \text{let } x_2 = 7 \text{ in } x_1 x_2, \text{let } x_3 = 7 \text{ in } x_1 x_3, \dots \} \end{aligned}$$

These notations can be generalized: Given some set A and some equivalence relation \approx , the name A/\approx refers to the set of equivalence classes, and for each element $a \in A$, the name $[a]_{\approx}$ refers to the set of all elements equivalent to a . Note that if $a_1 \approx a_2$, then $[a_1]_{\approx} = [a_2]_{\approx}$: they end up being *two different names for the same equivalence class*. Often if the equivalence class is obvious, then the subscript is dropped for succinctness, i.e., $[a_1] = [a_2]$

Equivalence classes give us a relatively clean way of brushing all this variable-renaming nonsense under the rug. The strategy essentially boils down to this: Rather than take the set TERM as the abstract syntax of our language, *take the alpha equivalence classes* TERM/\sim_a *to be the abstract syntax*.

Above, Prop. 8 showed that substitution respects alpha-equivalence classes. One consequence of this fact is that we can *define a substitution function over equivalence classes* $[t]_{\sim_a}$!¹²

$$\begin{aligned} \{\cdot/\cdot\} : (\text{TERM}/\sim_a) \times \text{VAR} \times (\text{TERM}/\sim_a) &\rightarrow (\text{TERM}/\sim_a) \\ \{[t]_{\sim_a}/x\}[t_0]_{\sim_a} &= [[t/x]t_0]_{\sim_a}. \end{aligned}$$

Substitution on alpha-equivalence classes is defined in terms of substitution on plain ole’ terms. Can you see why Prop. 8 makes this a well-defined function? At this point, we can redefine the big-step semantics to operate on equivalence classes of terms rather than particular Terms.

Well, this seems like a big mess. How did it help us at all? The actual answer is that it hasn’t yet. Our goal is to sweep some details under the carpet, but so far we’ve been dwelling on the details themselves. It’s now time to do some sweeping.

In the literature, you will typically not see this level of detailed explanation. Instead, what happens is this: The BNF for TERM is given, followed by a statement of the form “we consider terms *up to* alpha-equivalence” or “we consider terms *modulo* alpha-equivalence.” These magic words evoke much of the incantation above. From this point, any reference to TERM is really talking about TERM/\sim_a . Furthermore, any reference to a metavariable like t is really a reference to the alpha-equivalence class $[t]_{\sim_a}$.

Here’s where all this work shines through. First, a textual convention is established regarding metavariables. Remember that t really stands for $[t]_{\sim_a}$, and that there are many possible t ’s that one can consider as a name for the class (any TERM that is alpha-equivalent). The following convention is taken: “If t_1, t_2, \dots, t_n appear in a certain mathematical context, then in these terms all bound variables are chosen to be different from the free variables.” Since we are operating on alpha-equivalence classes, it’s possible to make this assumption because we can always find equivalence class representatives t_i that together satisfy this convention.

This approach, when treated informally, is called Barendregt’s Variable Convention, after the logician Henk Barendregt who introduced it (or just “The Variable Convention.”). In fact, most texts don’t bother mentioning that they assume the variable convention: instead they just note that they treat terms up to alpha-equivalence, and assume the variable convention without comment.

In practice, the variable convention means that whenever you write down terms or expressions that contain metavariables that stand for terms, you can safely assume that the bound variables of any subterm are distinct from the free variables of any subterm (because if they did not, then you could always alpha-

¹²We use curly braces for substitution here just so the definition looks clear. In practice, one would really use square brackets to talk about this substitution function.

convert to make that so). Under this convention, the definition of substitution becomes very short:

$$\begin{aligned}
 [t/x] : \text{TERM} &\rightarrow \text{TERM} \\
 [t/x]\text{true} &= \text{true} \\
 [t/x]\text{false} &= \text{false} \\
 [t/x](\text{if } t_1 \text{ then } t_2 \text{ else } t_3) &= \text{if } [t/x]t_1 \text{ then } [t/x]t_2 \text{ else } [t/x]t_3 \\
 [t/x]x &= t \\
 [t/x]x_0 &= x_0 \text{ if } x \neq x_0 \\
 [t/x](\text{let } x_0 = t_1 \text{ in } t_2) &= \text{let } x_0 = [t/x]t_1 \text{ in } [t/x]t_2
 \end{aligned}$$

The variable convention only comes into play in the last rule, where you silently assume that $x_0 \notin \{x\} \cup FV(t)$. Thus the renaming from our previous definition is implicitly applied if needed, and the relevant rule from substitution over actual TERMS applies.

This convention then carries over to definition of the big-step semantics for the language, as well as any other function or relation definitions. The result is a very concise definition on paper, but it's important to know what underlying hidden details are implied, especially that you are not operating on TERMS anymore. When exploiting this approach, it is *extremely important* that you prove that every step of evaluation respects α -equivalence. Once you show that that is true, you can abstract away.

An interesting side-effect of lifting to α -equivalence classes is that in an actual implementation of the language, you no longer care *which specific strategy* you use in practice to choose names. This may sound a little weird, because we went through the trouble of creating a *specific* strategy, where we ordered variables, etc. But the whole reason for doing that is just to show that there is *at least one* strategy for picking names that "works." Then, by abstracting to equivalence classes, we are more or less saying "we don't *really* care which concrete strategy you use, so long as it respects alpha equivalence, oh and by the way here is one specific one, just to show you that it's possible to pull it off. Do whatever you like! KThxBye!"

This is a common theme in building mathematical models as they get more sophisticated. We set down some properties that we want to be true of a system. Then we build one proof-of-concept variant to show that it's even possible (most language implementations could be viewed as such proofs-of-concept). Then we say that we only care about some class of properties of this system, and ignore others, so you can build a different implementation that works differently (so long as the differences fall in the category of "things-we-don't-care-about") but does the right thing as far as the specification is concerned. This is very analogous to Standards Documents and committees, where many implementations of a protocol, language, or piece of software (ostensibly) meet the criteria that are required, while otherwise doing whatever crazy town thing they like.