

CSE341, Fall 2011, Lecture 25 Summary

Standard Disclaimer: This lecture summary is not necessarily a complete substitute for attending class, reading the associated code, etc. It is designed to be a useful resource for students who attended class and are later reviewing the material.

Types for Objects (in the Next Lecture)

We previously studied static types for functional programs, in particular ML's type system. ML uses its type system to prevent errors like treating a number as a function. A key source of expressiveness in ML's type system (not rejecting too many programs that do nothing wrong and programmers are likely to write) is *parametric polymorphism*, also known as *generics*.

So we should also study static types for object-oriented programs, such as those found in Java. If everything is an object (which is less true in Java than in Ruby), then the main thing we would want our type system to prevent is “method missing” errors, i.e., sending a message to an object that has no method for that message. If objects have fields accessible from outside the object (e.g., in Java), then we also want to prevent “field missing” errors. There are other possible errors as well, like calling a method with the wrong number of arguments.

While languages like Java and C# have generics these days, the source of type-system expressiveness most fundamental to object-oriented style is *subtype polymorphism*, also known as *subtyping*. ML does not have subtyping, though this decision is really one of language design (it would complicate type inference, for example).

Our plan for this lecture and the next one is to:

- Study subtyping
- Compare subtyping and generics, determining which idioms are best supported by each
- Combine subtyping and generics, showing that the result is even more useful than the sum of the two techniques

Subtyping for Records and Functions

It would be natural to study subtyping using Java since it is a well-known object-oriented language with a type system that has subtyping. But it is also fairly complicated, using classes and interfaces for types that describe objects with methods, overriding, static overloading, etc. While these features have pluses and minuses, they can complicate the fundamental ideas that underlie how subtyping should work.

So this lecture studies subtyping using only records (like in ML, things with named fields holding contents — basically objects with public fields, no methods, and no class names) and functions (like in ML or Racket). This will let us see how subtyping should — and should not — work, so that then the next lecture can then apply the results to the more complicated setting of a class-based object-oriented language.

This approach has the disadvantage that we cannot use any of the language we have studied: ML does not have subtyping and record fields are immutable, Racket and Ruby are dynamically typed, and Java is too complicated for our starting point. So we are going to make up a language with just records, functions, variables, numbers, strings, etc. and explain the meaning of expressions and types as we go.

A Made-Up Language of Records

To study the basic ideas behind subtyping, we will use records with mutable fields, as well as functions and other expressions. Our syntax will be a mix of ML and Java that keeps examples short and, hopefully, clear. For records, we will have expressions for making records, getting a field, and setting a field as follows:

- In the expression $\{f1=e1, f2=e2, \dots, fn=en\}$, each f_i is a field name and each e_i is an expression. The semantics is to evaluate each e_i to a value v_i and the result is the record value $\{f1=v1, f2=v2, \dots, fn=vn\}$. So a record value is just a collection of fields, where each field has a name and a contents.
- For the expression $e.f$, we evaluate e to a value v . If v is a record with an f field, then the result is the contents of the f field. Our type system will ensure v has an f field.
- For the expression $e1.f = e2$, we evaluate $e1$ and $e2$ to values $v1$ and $v2$. If $v1$ is a record with an f field, then we update the f field to have $v2$ for its contents. Our type system will ensure $v1$ has an f field. Like in Java, we will choose to have the result of $e1.f = e2$ be $v2$, though usually we do not use the result of a field-update.

Now we need a type system, with a form of types for records and typing rules for each of our expressions. Like in ML, let's write record types as $\{f1:t1, f2:t2, \dots, fn:tn\}$. For example, $\{x : \text{real}, y : \text{real}\}$ would describe records with two fields named x and y that hold contents of type real . And $\{\text{foo} : \{x : \text{real}, y : \text{real}\}, \text{bar} : \text{string}, \text{baz} : \text{string}\}$ would describe a record with three fields where the foo field holds a (nested) record of type $\{x : \text{real}, y : \text{real}\}$. We then type-check expressions as follows:

- If $e1$ has type $t1$, $e2$ has type $t2$, ..., en has type tn , then $\{f1=e1, f2=e2, \dots, fn=en\}$ has type $\{f1:t1, f2:t2, \dots, fn:tn\}$.
- If e has a record type containing $f : t$, then $e.f$ has type t (else $e.f$ does not type-check).
- If $e1$ has a record type containing $f : t$ and $e2$ has type t , then $e1.f = e2$ has type t (else $e1.f = e2$ does not type-check).

Assuming the "regular" typing rules for other expressions like variables, functions, arithmetic, and function calls, an example like this will type-check as we would expect:

```
fun distToOrigin (p:{x:real,y:real}) =
  Math.sqrt(p.x*p.x + p.y*p.y)

val pythag : {x:real,y:real} = {x=3.0, y=4.0}
val five : real = distToOrigin(pythag)
```

In particular, the function `distToOrigin` has type $\{x : \text{real}, y : \text{real}\} \rightarrow \text{real}$, where we write function types with the same syntax as in ML.

This type system does what it is intended to do: No program that type-checks would, when evaluated, try to look up a field in a record that does not have that field.

Now Add Subtyping

With our typing rules so far, this program would not type-check:

```
fun distToOrigin (p:{x:real,y:real}) =
  Math.sqrt(p.x*p.x + p.y*p.y)
val c : {x:real,y:real,color:string} = {x=3.0, y=4.0, color="green"}
val five : real = distToOrigin(c)
```

In the call `distToOrigin(c)`, the type of the argument is $\{x:\text{real},y:\text{real},\text{color}:\text{string}\}$ and the type the function expects is $\{x:\text{real},y:\text{real}\}$, breaking the typing rule that functions must be called with the

type of argument they expect. Yet the program above is safe: running it would not lead to accessing a field that does not exist.

A natural idea is to make our type system more lenient as follows: If some expression has a record type $\{f_1:t_1, \dots, f_n:t_n\}$, then let the expression *also* have a type where some of the fields are removed. Then our example will type-check: Since the expression c has type $\{x:\text{real}, y:\text{real}, \text{color}:\text{string}\}$, it can also have type $\{x:\text{real}, y:\text{real}\}$, which allows the call to type-check. Notice we could also use c as an argument to a function of type $\{\text{color}:\text{string}\} \rightarrow \text{int}$, for example.

Letting an expression that has one type also have another type that has less information is the idea of *subtyping*. (It may seem backwards that the *subtype* has *more* information, but that is how it works. A less-backwards way of thinking about it is that there are “fewer” values of the subtype than of the supertype because values of the subtype have more obligations, e.g., having more fields.)

We will now add subtyping to our made-up language, in a way that will not require us to change any of our existing typing rules. For example, we will leave the function-call rule the same, still requiring that the type of the actual argument *equal* the type of the function parameter in the function definition. To do this, we will add two things to our type system:

- The idea of one type being a subtype of another: We will write $t_1 <: t_2$ to mean t_1 is a subtype of t_2 .
- One and only new typing rule: If e has type t_1 and $t_1 <: t_2$, then e (also) has type t_2 .

So now we just need to give rules for $t_1 <: t_2$, i.e., when is one type a subtype of another.

Subtyping is Not a Matter of Opinion

A common misconception is that if we are defining our own language, then we can make the typing and subtyping rules whatever we want. That is only true if we forget that our type system is allegedly preventing something from happening when programs run. If our goal is (still) to prevent field-missing errors, then we cannot add any subtyping rules that would cause us to stop meeting our goal.

For subtyping, the key guiding principle is *substitutability*: If we allow $t_1 <: t_2$, then any value of type t_1 must be able to be used in every way a t_2 can be. For records, that means t_1 should have all the fields that t_2 has and with the same types.

Some Good Subtyping Rules

Without further ado, we can now give four rules that we can add to our language to accept more programs without breaking the type system. The first two are specific to records and the next two, while perhaps seeming unnecessary, do no harm and are common in any language with subtyping because they combine well with other rules:

- “Width” subtyping: A supertype can have a subset of fields with the same types, i.e., a subtype can have “extra” fields
- “Permutation” subtyping: A supertype can have the same set of fields with the same types in a different order.
- Transitivity: If $t_1 <: t_2$ and $t_2 <: t_3$, then $t_1 <: t_3$.
- Reflexivity: Every type is a subtype of itself: $t <: t$.

Notice that width subtyping lets us forget fields, permutation subtyping lets us reorder fields (e.g., so we can pass a $\{x:\text{real}, y:\text{real}\}$ in place of a $\{y:\text{real}, x:\text{real}\}$) and transitivity with those rules lets us do both (e.g., so we can pass a $\{x:\text{real}, \text{foo}:\text{string}, y:\text{real}\}$ in place of a $\{y:\text{real}, x:\text{real}\}$).