

Lecture 2: Object Semantics

These notes cover two lectures on the topic of *object semantics*. The material is chosen with three purposes in mind:

- to help you become familiar with the basic runtime mechanism common to all object-oriented languages (but with a particular focus on Java): variables, object references, assignments, mutability, and so on.
- to introduce two diagrammatic notations, *object diagrams* for describing ‘snapshots’ (that is, particular configurations of objects in the heap), and—far more useful—*object models*, for describing sets of snapshots.
- to make you aware, in passing, of some tricky issues that we’ll return to later in more detail, and which turn out to be of fundamental importance: equality, rep invariants and exposure, and subtype polymorphism.

When you’ve completed this material, you should have a solid grasp of what happens when Java code executes, so that you can predict what some code will do without running it. You should be able to read object models, but it will take some practice before you’re confident writing them.

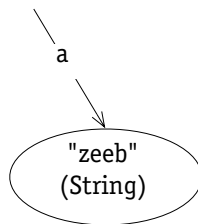
2.1 Variables, References and Objects

Some types of objects can be created with *literals*. What happens when you run this?

```
2.1.1 String a = "zeeb";  
2.1.2 System.out.println (a);
```

It prints `zeeb`. Statement 2.1.1 is a declaration (of the variable `a`) and an assignment (to `a`) in one. The expression `"zeeb"` is called a string literal. Statement 2.1.2 is a procedure call that prints a representation of `a` to standard out; don’t worry for now about its details.

We can draw the result of the first statement as an *object diagram*, showing that `a` is a *reference* to an object of type `String`:

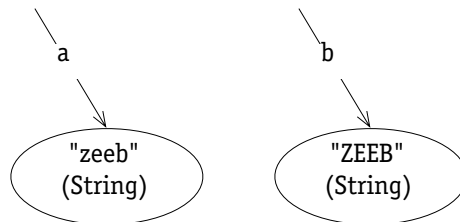


That `a` is a reference to the string (and not the string itself, or a slot that holds it) is important to understand, but actually not very significant in this case.

What happens when you run this?

```
2.1.3 String a = "zeeb";  
2.1.4 String b = a.toUpperCase ();  
2.1.5 System.out.println (b);
```

It prints ZEEB. Statement 2.1.4 is a call to the method `toUpperCase`. The method has one argument, `a`, which sits in a special position: it's called the 'receiver'. The call results in the creation of a fresh string that is then bound to the variable `b`. Here's the object diagram:



What does this do?

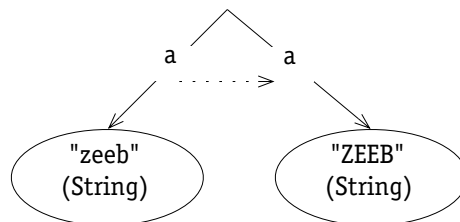
```
2.1.6 String a = "zeeb";  
2.1.7 a.toUpperCase ();  
2.1.8 System.out.println (a);
```

It prints `zeeb`. Statement 2.1.7 creates a fresh string that gets thrown away since it is bound to no variable. The string object referenced by `a` is not changed: strings are immutable.

What about this?

```
2.1.9 String a = "zeeb";  
2.1.10 a = a.toUpperCase ();  
2.1.11 System.out.println (a);
```

Again, no object changes. But `a` is made to refer to the new object by the assignment 2.1.10, so the result is ZEEB.

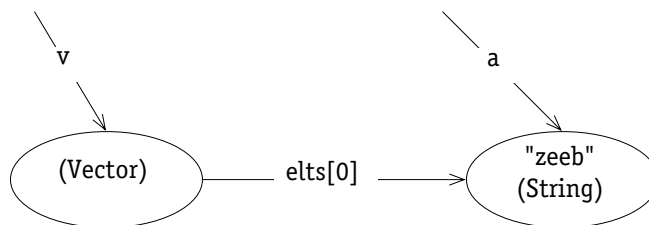


2.2 Aliasing, Mutability and Reference Equality

Java provides a variety of collections as part of its standard library. A vector is like an array, but it can grow and shrink dynamically. An example of using a vector:

```
2.2.1 Vector v = new Vector ();
2.2.2 String a = "zeeb";
2.2.3 v.add (a);
2.2.4 System.out.println (v.lastElement ());
```

This prints `zeeb`, because statement 2.2.3 inserts (a reference to) the string (referred to by) `a` into the vector (referred to by) `v`:

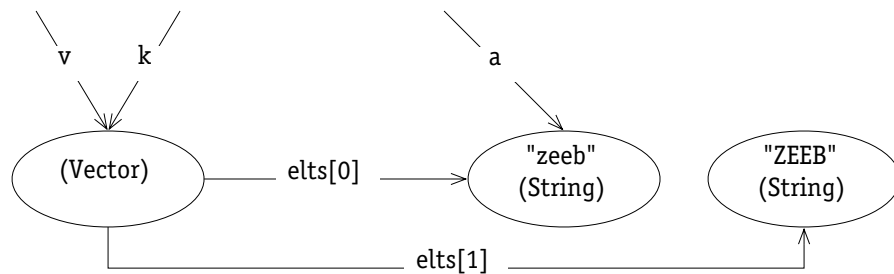


and statement 2.2.4 takes this element and prints it. The method `lastElement` is like the string method `toUpperCase`: it takes a single argument, the receiver, has no effect on it, and returns a reference to an object (in this case, the last element of the vector `v`). The method `add`, on the other hand, takes two arguments: the first is the receiver `v`, and the second is the string argument `a`. (Sometimes we'll say that such a method takes only *one* argument, ignoring the receiver.) Unlike `toUpperCase`, the method `add` does not return an object, but changes or *mutates* its receiver object. Vectors, unlike strings, can change, and are thus said to be *mutable*.

To see this, consider:

```
2.2.5 Vector v = new Vector ();
2.2.6 Vector k = v;
2.2.7 String a = "zeeb";
2.2.8 v.add (a);
2.2.9 k.add (a.toUpperCase());
2.2.10 System.out.println (v.lastElement ());
```

This prints `ZEEB`, since after statement 2.2.6 the two variables `k` and `v` are names for the same vector object: they are said to be *aliases*. The calls to `add` mutate the one vector object, first adding the object for the lower case string, then the upper case string. The changes are visible through both names, so in statements 2.2.8 to 2.2.10 we could actually permute the names `k` and `v` without any change in behaviour.



Aliasing is pervasive in languages like Java, and very useful. But it adds a lot of complexity. For one thing, it breaks the rule that a statement ‘affects only the variables it mentions’. Just because `v` isn’t mentioned in statement 2.2.9 doesn’t mean that it won’t affect the result of statement 2.2.10 which mentions only `v` and not `k`.

How can we observe the aliasing more directly? By testing equality:

```

2.2.11 Vector v = new Vector ();
2.2.12 Vector k = v;
2.2.13 if (v == k)
2.2.14     System.out.println ("same");
2.2.15 else
2.2.16     System.out.println ("different");
2.2.17
  
```

which results in `same` being printed. The built-in `==` test tells you whether two references are for the same object, so it’s often called a test of ‘reference equality’.

What does this do?

```

2.2.18 Vector v = new Vector ();
2.2.19 Vector k = new Vector ();
2.2.20 if (v == k)
2.2.21     System.out.println ("same");
2.2.22 else
2.2.23     System.out.println ("different");
  
```

It prints `different`, because `v` and `k` are distinct objects. It’s a fundamental property of constructors that the objects they return really are fresh. In fact, the garbage collector can recycle an object, but only if there is no reference to it still around. This ensures that even if objects are recycled, we can never tell. Here are the object diagrams for the two cases:

A puzzle: what does this do?

```

2.2.24 String a = "zeeb";
2.2.25 String b = "zeeb";
  
```

```

2.2.26  if (a == b)
2.2.27      System.out.println ("same");
2.2.28  else
2.2.29      System.out.println ("different");

```

Strangely, this prints same, because the Java virtual machine automatically ‘interns’ string literals: if it can tell that two string literals have the same sequence of characters, it only allocates one object. You’d be right to think this is a bit confusing; it’s a performance optimization.

In fact, it’s very bad form to test reference equality of immutable objects, unless you’re doing something subtle with memory management. Arguably it’s a design defect of Java that you can even observe whether two immutable objects are the same or not. So how should you compare two immutable objects? With an equals method.

The String class provides a method equals that tells you whether two strings contain the same sequence of characters or not. This code

```

2.2.30  String a = "zeeb";
2.2.31  String b = a.toUpperCase ();
2.2.32  if (b.equals ("ZEEB"))
2.2.33      System.out.println ("same characters");
2.2.34  else
2.2.35      System.out.println ("different characters");

```

prints same characters.

When we study inheritance, you’ll learn that every class automatically inherits an equals method, so you might think you don’t need to write one. But it’s almost never what you want, so whenever you design a class, one of the first things you’ll need to figure out is when two objects of the class should be considered equal to one another.

Break for questions:

- Would you expect that generally $x == y$ implies $x.equals(y)$? Yes, it should. Because the equals method can be user-defined, just like any other method, you could make it behave in any way you wanted. On a mutable type, it might even mutate the object! But that would be disastrous: there’s a generic contract that clients expect equals to obey. More on this later.
- Why would a language have immutable types? Because aliasing is complicated, and when you use immutable types, the issue doesn’t arise. Also, code built with immutable types can sometimes be more efficient.
- So if immutable types are so much simpler, why have mutable types? Because mutation gives a very useful form of modularity: it allows you to make local changes to a structure. And mutation is often a natural way to model entities in the real world: a transaction on a bank account changes it; it doesn’t produce a new bank account.

(Bob Harper of Carnegie Mellon, one of the designers of the programming language ML, pointed out to me that people have a misconception that mutation is all about performance—because often it’s more efficient to mutate a large structure in one place than to replace the entire structure—but it’s actually about expressiveness.)

2.3 Null References

What does this do?

```
2.3.1 String a = null;
2.3.2 System.out.println (a);
```

It prints null. The keyword null denotes a value that can be taken on by an object reference. It means that the reference does not in fact refer to any object. There is no null object!

But note that this code

```
2.3.3 String a = null;
2.3.4 String b = a.toUpperCase ();
2.3.5 System.out.println (b);
```

behaves quite differently. It throws a `NullPointerException` on line 2.3.4. We’ll learn about exceptions later, but for now, all you need to understand is that line 2.3.4 failed, when the expression `a.toUpperCase()` was evaluated.

What’s the difference? The receiver to a method call can never be null, because it identifies the object that ‘receives’ the call—and that has to be some object. So `a.toUpperCase()` fails when `a` is null. But in the previous example, `System.out.println(a)` is OK when `a` is null, since reference is an argument, and this is a special kind of method (really just a plain old procedure) that doesn’t have a receiver.

You can write a method that tests whether an argument is null and does something appropriate. Dereferencing null is a common programming mistake in Java. To avoid it, you can check whether a reference is null before you attempt to call a method:

```
2.3.6 if (arg == null)
2.3.7     System.out.println ("error");
2.3.8 else {
2.3.9     String x = arg.toUpperCase ();
2.3.10 ...
```

This is how `System.out.println` works; it does something like this:

```
2.3.11 if (arg == null)
2.3.12     print ("null");
```

```
2.3.13 else {
2.3.14     ...
2.3.15 }
2.3.16 ...
```

In general, rather than catching nulls and treating them specially, it's better to avoid creating null references in the first place. You'll learn about that when we discuss representation invariants. Sometimes you can't avoid it, and then it's important to document where the null references may occur. That's one reason specifications are important: they can spare you runtime errors and unnecessary checks.

Questions:

- In general, would you expect `a.equals (b)` to be substitutable for `b.equals (a)`? No, because when `a` is null and `b` is not, the first will throw an exception, and the second will (usually) return false.
- OK, smarty pants, leave nulls alone. Would you then expect `a.equals (b)` and `b.equals (a)` to have the same effect? Yes, you would. In fact, this property of the equals method—called *symmetry*—is demanded by Java's 'object contract'. We'll see later when we study equality in depth what other properties are required, and how easy it is to mess up and write an equals method that does *not* have these properties.

2.4 User-defined Classes & Fields

Let's make some objects of our own:

```
2.4.1 class Trans {
2.4.2     int amount;
2.4.3     Date date;
2.4.4 }
2.4.5
```

This code declares a *class*, a kind of template for making objects. These objects are going to represent transactions in a banking system. Each object has an integer amount (which may be negative for a withdrawal), and a date/time stamp to mark the moment at which the transaction occurred. We'll use these transaction objects within an account object, so we don't need to associate an account explicitly with each transaction.

The class declares two *fields* or *instance variables*, `amount` and `date`. Each object of the class will contain two references, one to an integer and one to a date. The type `Date` is a class from the Java library; it's predefined like `String` (but not part of the language definition the way `String` is).

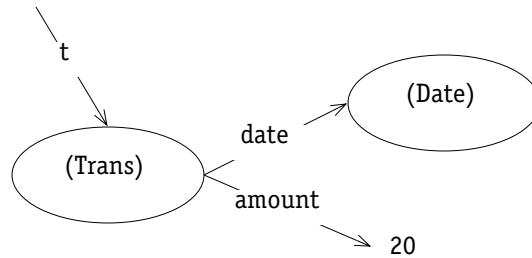
The type `int` is a rather strange beast. It's not a class at all, but a *primitive type*. Variables or fields of type *int* don't hold references to integer objects; they hold the integers them-

selves. You may think it a bit jarring that an object-oriented language has this rather unobject-oriented notion in it (and many people share your opinion). Sometimes we'll actually need an integer that's an object, and in that case we can use the class `Integer` from the Java library. How do you get from an `int` to an `Integer` and back? Well, that's a little obscure, and that's one reason people don't like this design.

If we run this code:

```
2.4.6   Trans t = new Trans ();
2.4.7   t.amount = 20;
2.4.8   t.date = new Date ();
```

a fresh object gets created, and its fields are set, resulting in this configuration:



The expression on the right-hand side of statement 2.4.6 is a call to a *constructor*: it makes a new object with default values for the fields. In this case, it creates a `Trans` object with zero for the amount and null for date. The statement 2.4.7 is called a *setter*: it sets the value of the field `amount` of the object referred to by `t`. Statement 2.4.8 has another constructor call on the right; creating a new `Date` by default creates a `Date` object representing the moment at which the the object is itself created. But it's also a setter: it sets the `date` field of `t` to point to this new date.

2.5 User-defined Constructors

So we've succeeded in making a `Trans` object representing a deposit of twenty dollars at this moment in time. The way we did it – creating an uninitialized object and then setting its fields – is not a good one, however. We'll want our transactions to be well-formed; for example we won't want to have transactions that don't have dates. And perhaps we'll want every transaction to have a non-zero amount. Later, we'll study these kinds of *invariants* in much more depth.

For now, just observe that immediately after statement 2.4.6 we have a transaction object that is *not* well formed. When an object is created, its fields are initialized to default values: null for object references, and zero for integers. So `t.amount` will be zero, and `t.date` will be null. These default values are rarely what you want; after all, which values

make sense will depend on the problem we're trying to solve. In this case, you'd have to know something about banking to know that a transaction of zero dollars is ill-formed.

Is it a big deal that there's a bogus transaction object hanging around between statements 2.4.6 and 2.4.7? Yes, it is, and here's why. We'd like the responsibility for ensuring that objects of the class `Trans` are well formed to be handled entirely within the `Trans` class. In our program, you need to check not only the code of the class, but also the code that uses the class. At this scale, it's not a disaster. But in a much larger program, you need as much *modularity* as you can get, confining tricky aspects of the program as much as possible to small areas of the code.

To solve this problem, we declare our own constructor:

```
2.5.1    class Trans {
2.5.2        int amount;
2.5.3        Date date;
2.5.4        Trans (int a, Date d) {amount = a; date = d;}
2.5.5    }
```

This constructor takes an amount and a date as arguments, and creates a transaction object with that amount and date. If I'd wanted to ensure that the amount of a transaction is non-zero, I could have added a check that threw an exception if the amount was zero; we'll see how to do that later.

A peculiar, but useful, property of constructors is that having defined our own constructor, the default constructor – the one that just initializes each field to default values – becomes no longer available. So statement 2.4.6 will no longer compile. Instead, we can write

```
2.5.6    Trans t = new Trans (20, new Date ());
```

and that one line of code will have the effect that lines 2.4.6 to 2.4.8 had previously.

(We haven't fully solved the problem of modularizing the invariant of `Trans`, by the way. You can still mess up a `Trans` object by setting its date field to null from outside, for example. The first step to prevent this is to make use of Java's accessibility mechanisms: we can make the fields *private* so they can't be read or written from outside the class. In fact, this turns out not to be enough, as we'll learn when we study data abstraction.)

2.6 Method Call

Here's a class representing a bank account, with fields for the name of the account, a vector of transactions, and the account balance:

```
2.6.1    class Account {
2.6.2        String name;
```

```

2.6.3   Vector transv;
2.6.4   int balance;
2.6.5   Account (String n) {
2.6.6       transv = new Vector ();
2.6.7       balance = 0;
2.6.8       name = n;
2.6.9       }
2.6.10  boolean checkTrans (Trans t) {
2.6.11      return (balance + t.amount = 0);
2.6.12  }
2.6.13  void post (Trans t) {
2.6.14      transv.add (t);
2.6.15      balance += t.amount;
2.6.16  }
2.6.17  }

```

There are two methods: to post a transaction, and to check whether a transaction is legal (or would result in a negative balance if posted). Let's investigate what happens when such methods are called.

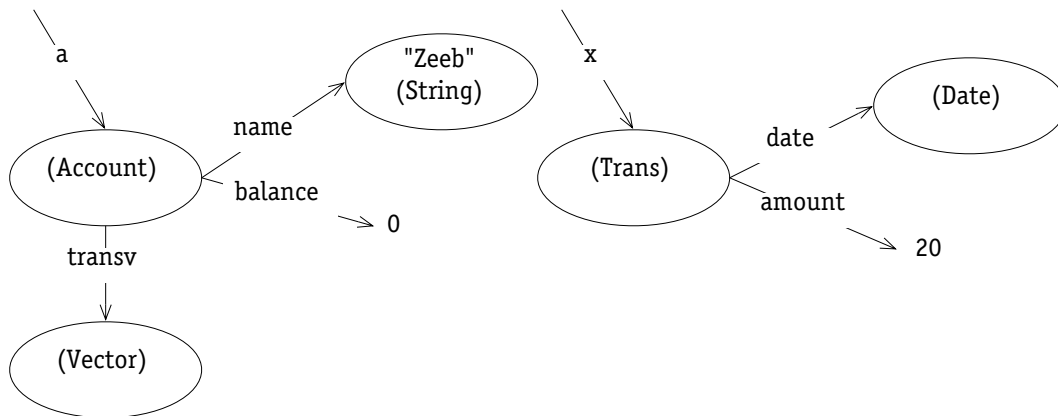
Consider this code fragment:

```

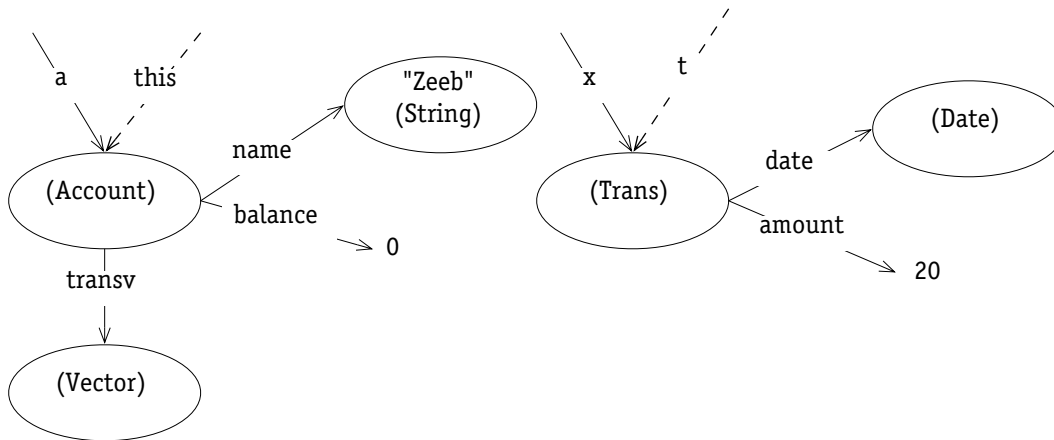
2.6.18  Account a = new Account ("Zeeb");
2.6.19  Trans x = new Trans (20, new Date ());
2.6.20  if (a.checkTrans(x)) {
2.6.21      a.post (x);
2.6.22  }

```

Let's look at the effects of the method calls. Just before the call on line 2.6.20, we're in this configuration:



The first step in the call is to bind the arguments: the argument `t` is made, within the method, to reference the object called `x` outside the method, and the special argument `this`, representing the receiver, is made to reference the object `a`:



Several things to note here:

- Java is lexically scoped (like Scheme), so the variables `a` and `x` are not available inside the method call.
- Like any procedure call, the argument variables are allocated on the stack, so they are only available for the duration of the method call. The variable `this` may of course be available outside, but will refer to a different object.
- The formal arguments are bound as fresh references to the objects referenced by the expressions passed as actual arguments, so that they ‘share’ objects. The outside variable `x` and the formal argument `t`, for example, share a `Trans` object. For this reason, the term *call by sharing* is sometimes used for this form of parameter passing. Alternatively, you can view it as *call by value*, in which *references* are passed.

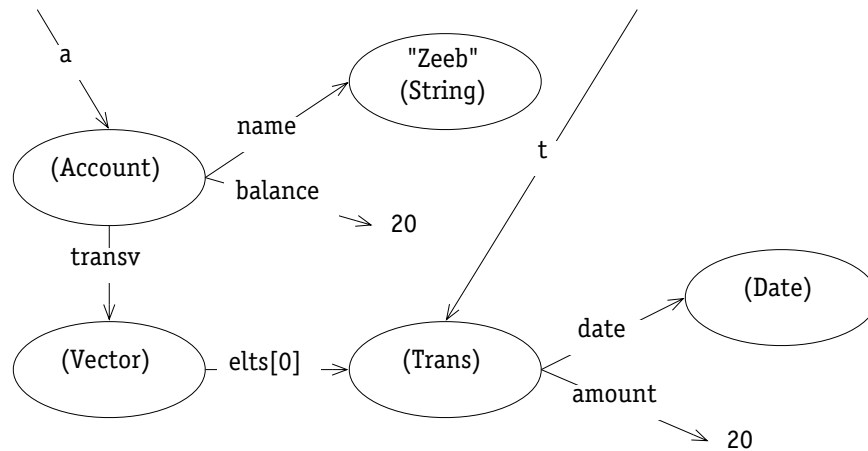
To execute the body of `checkTrans`, we now execute the statement

```
2.6.23 balance + t.amount = 0
```

in the context of the bindings we’ve just created. When the name of a field `f` appears within the code of its class unqualified by an object reference, it is taken to be short for `this.f`. So here `balance` is short for `this.balance`, which evaluates to the value of the `balance` field of the object called `a` outside. The expression `t.amount` gets the value of the `amount` field of the object called `x` outside. The result of the entire expression is `true`, and is returned as the value of the method call, and the body of the `if`-statement now executes.

The call on line 2.6.21 proceeds by exactly the same mechanism. This time, however, two objects are mutated: the `Account` object, whose `balance` field is incremented, and

the Vector object in its transv field, to which the transaction is added as a new element. The result of the call by sharing is that, after the call, the object a has changed:



You may be wondering what significance there is to the special ‘receiver’ argument. So far, it really has no significance; method call is really just procedure call, with some syntactic shorthands. The significance arises when we look at subclassing in our next lecture, where several methods can have the same name, and which gets called will depend on properties of the receiver argument alone.

2.7 Conclusion

Now we’ve seen all the key notions for how objects are manipulated in an object-oriented language. We’ve seen how they’re created; how references are bound to objects; and how their fields are set. We’ve mentioned the two fundamental kinds of quality – reference equality (tested with `==`) and object equality (tested with an `equals` method – about which we’ll have much more to say later. And we’ve looked at parameters are passed in method calls. These mechanisms comprise one important part of what it means to be ‘object oriented’.