

## 5 Errors and Exceptions

### 5.1 Kinds of Errors

Consider the kind of errors that can occur in programs, generally speaking. There are various errors possible, and it makes sense to classify them somehow.

- (1) Errors in syntax. For instance, writing `clas` instead of `class`.
- (2) Applying an operation that works only on some type of values on values of the wrong type. For instance, trying to add two Boolean values, or trying to divide strings.
- (3) Invoking a method `m` on an object that does not define such a method `m`, for instance, if `d` is an instance of `Drawing`, trying to invoke `d.playTune()` is an error.
- (4) Invoking an operation with arguments of the right types for which the operation is undefined. For instance, dividing by 0, or taking the tangent of  $\pi/2$ .
- (5) Problems out of a programmer's control, such as hardware failures. For instance, disk failures during disk IO, network failure during network IO, or the CPU melting. Running out of memory is also something that sometimes occurs.
- (6) "Logical" errors, where the program does not behave as expected, such as an implementation of addition that returns 3 when given inputs 1 and 1, or a `draw` method on a drawing that flips all the lines upside down.

These distinctions are not sharply defined; some errors may well be classified in different categories. But the categories are useful as a general sense of the kind of errors that arise.

Different languages check and deal with these errors differently. The first kind of errors, syntax errors, are usually caught before a program is executed — sometimes even caught by whatever programming interface you are using for your language.

The other kind of errors can be caught either before a program executes, when it is compiled (we say that these errors are handled *statically*), or when a program executes (we say these errors are handled *dynamically*). Take errors of the kind (2) and (3). When a language checks for those kind of errors dynamically, we say that the language supports *dynamic type checking*. Scheme supports dynamic type checking. Other languages actually can check for

type errors at compile time, that is, before programs execute. They support *static type checking*. Java supports static type checking.

There are advantages to checking for errors at compile time; in particular, you still have a chance to correct the problem while the code is in your hands. With dynamic type checking, errors may only show up after the code has been shipped and the piece of software is in the hands of the customer, making it more difficult and expensive to correct. Static error checking also has some disadvantages. In particular, there is no way to identify exactly all those programs that have errors at compile time.<sup>1</sup> Thus, the error checker needs to approximate, and it will approximate conservatively. Thus, there are programs that would not cause problems during execution that the error checker will reject.

Java takes care of type errors, errors of type (2) and (3) above, at compile time. The other kind of errors (4) and (5), however, cannot be reliably checked for at compile time. Here again, the limitation in the footnote above bites us here. Every language will have a different way to report those kind of errors. Generally, the error will abort execution and report a useful error message on the console. But in many languages, Java included, these errors can be dealt with within the program itself, and the execution need not actually abort. In other words, these errors can often be recovered from gracefully. Other errors such as (6) are not even caught by the language, since they rely on the notion of “what’s expected” which only exists in the designer’s or programmer’s head. I’ll talk about logical errors next lecture.

## 5.2 Exceptions

In medieval times, one way to report errors at run time was to have functions return a special value called various an *error value*, or an *error flag*. Basically, a special value that did not make sense as a return value for the function, and that could be used to indicate that execution failed. For instance, a function could return  $-1$  to indicate an error instead of the usual positive number it would return on a correct execution.

The problem with error values is that it is very easy to forget to check the result of a function call to see if an error occurred. And forgetting to check the result of a function call means that if there was indeed an error, the error value gets propagated in your code and will likely cause problems somewhere else (e.g., if you end up taking the square root of the value returned by the function, and that value was the error value  $-1$ ). Tracking these kinds of bugs is difficult.

This error-value technique is still in use nowadays, unfortunately. In Java, for instance, the *null* value is often used as an error marker — the *null* value is a special object that can be

---

<sup>1</sup>This is a deep limitation in what we can say about programs in general, which you will see in a good theory of computation course. It is a consequence of the so-called *undecidability of the halting problem for Turing machines*. Very roughly speaking, the limitation is that it is impossible to write a program that takes a program  $P$  as input and (without executing the program) answers correctly a question about how  $P$  will execute.

considered an instance of every class. The problem is that, as I mentioned, it is very easy to forget to check that a value is not *null* when returned from a function, and that *null* value can be propagated in the code, and cause problems when someone tries to invoke a method on *null*, which does not implement any method. I strongly suggest you avoid using *null* for error reporting — in fact you should probably avoid *null* altogether. (We'll see techniques for doing so later.)

The error-value technique has the consequence that if you don't check that an error occurred, that error is free to wreak havoc in the rest of your program. That's bad, because programmers are lazy and will forget to check for errors.

The opposite point of view is taken in the modern approach to dealing with run-time errors, *exceptions*. An exception is just an object in the system, that gets created when an error is encountered, and that propagates through the code until either it is handled, or aborts execution. Thus, if an exception is not dealt with, instead of letting your code continue and cause further problem, execution aborts completely. If you want your code to keep going — perhaps because you have another way to try to achieve your goal, such as trying to write on the network if writing to disk failed — that's when you have to put in some work.

Most modern languages have exceptions, and they all essentially work in the same way. Let's look at how Java implements exceptions, since it is actually fairly representative (as well as useful for us). There is actually a whole hierarchy of classes in Java implementing exceptions. This hierarchy lets us distinguish the kind of exceptions that can occur. Here is a partial class hierarchy of exceptions:

```
Throwable
|
+-- Error
|   |
|   +-- OutOfMemoryError
|   +-- AssertionError
|
+-- Exception
    |
    +-- IOException
    +-- InterruptedException
    +-- RuntimeException
        |
        +-- ArithmeticException
        +-- ClassCastException
        +-- NullPointerException
        +-- IllegalArgumentException
        +-- NumberFormatException
```

The class `Throwable` is the most general kind of exception that every other exception

subclasses. The `Error` class roughly represent the show-stoppers, that lead to aborting execution in almost all cases. The `Exception` class capture more “benign” forms of errors. These include `IOExceptions`, representing exceptions due to failure of IO (disk failure, network failure, and so on), while `RuntimeExceptions` represent exceptions such as dividing by 0 (an `ArithmeticException`), casting an object to an unacceptable class (a `ClassCastException`), invoking a method on a null object (a `NullPointerException`). The `IllegalArgumentException` is a general exception to represent passing a wrong value to a method.

(Note that you can also create new kinds of exceptions by subclassing, generally by subclassing the `Exception` class. We’ll see subclassing more extensively in the coming weeks.)

Many exceptions are created automatically by the system when an error is encountered. You can also cause an exception yourself in the code. This is called *throwing* an exception:

```
throw new IllegalArgumentException("Oops, something bad happened");
```

To a first approximation, throwing an exception aborts execution of the program. Here is a simple, possibly naive, but still accurate model that explains how exceptions affect code execution. Intuitively, you can think of the following two rules that apply to all programs:

- (1) Every method returns either a value as specified by its signature, or an exception. The statement `return` returns a value from the method, while the statement `throw` returns an exception from the method.
- (2) Every method call is implicitly wrapped by code that checks if the method called returned an exception—if so, it returns the exception immediately, using a `throw`, otherwise it just continues execution with the value returned by the method.

In this way, exceptions propagate all the way back to the `main` method. If the `main` method returns an exception, then that exception is reported to the user.

In fact, the fact that exceptions can be returned from methods sometimes has to appear in the signature. Java distinguishes between checked and unchecked exceptions. Unchecked exceptions (including `Errors` and `RuntimeExceptions`) are exceptions that can occur essentially at any point during execution. Checked exceptions (including `IOExceptions`) and all exceptions that you will create) are exceptions that can occur only when specific methods throw them. You must annotate every method that can throw a checked exception (either by throwing it directly or by invoking a method that can throw that exception):

```
public void someMethod () throws SomeException {  
    // some code that can throw the SomeException exception  
}
```

This is all good and well, but how can we deal with such exceptions gracefully? After all, if exceptions were just errors that cannot be dealt with, we could just say that an exception

aborts execution, and be done with it. The idea is that if a method call appear in the body of a `try` block, and if that call returns an exception, that exception is not returned immediately, but instead it is handled `catch` clause associated with the `try` block.

Suppose `MyException` is a new exception you have defined and suppose you wanted to intercept this exception if it is thrown by the method `someMethod()` on object `obj`:

```
try {
    obj.someMethod();
} catch (MyException e) {
    // some code to deal with the exception
}
```

This reads: if a `MyException` is returned by `obj.someMethod()`, then instead of having that exception returned immediately, intercept it, name it `e`, and continue execution of the current method with the code in the `catch` clause. (The reason why we may be interested in `e` is that `e` is an object that actually implements some useful methods such as `getMessage()` which returns a string representing the message associated with the exception.) This lets you gracefully recover from an exception, by intercepting it and dealing with it instead of letting it bubble up until it aborts the entire program with an error.

Note that a `catch` clause catching an exception `SomeException` will in fact catch all exceptions that are instances of subclasses of `SomeException`. This lets you intercept, for instance, any possible exception in a single `catch` clause:

```
try {
    // some code doing something interesting
} catch (Throwable e) {
    // deal with the exception
}
```

This works because every exception ultimately subclasses `Throwable`. This kind of code is useful in testing code to nicely format the exception and report useful information.

It is also possible to catch multiple exceptions and dealing with them differently:

```
try {
    // some code doing something interesting
} catch (SomeExceptionClass e) {
    // deal with this kind of exception
} catch (SomeOtherExceptionClass e) {
    // deal with this other kind of exception
}
```

The order of the `catch` clauses is relevant — they are tried in order, and the first clause where the current exception is in a subclass of the specified exception will be chosen.