

programming has become an important tool for the challenges of our time. As you, a Java developer, can use it to your advantage.

The recent interest in functional programming started as a response to the pervasiveness of *concurrency* as a way of scaling horizontally, through Multithreaded programming (see, e.g., [Goetz2006]) is difficult to do. Most developers are good at it. As we'll see, functional programming offers help for writing robust, concurrent software.

An example of the greater need for horizontal scalability is the growth of data sets requiring management and analysis, the so-called *big data* trend. These are data sets that are too large for traditional database management systems. The

Because this is a short introduction and because it is difficult to represent traditional concepts in Java, there will be several topics that I won't discuss, although I have added glossary entries, for completeness. These topics include *lazy evaluation*, *partial application*, and *comprehensions*. I'll briefly discuss several other topics, such as *combinators*, *laziness*, and *monads*, to give you a taste of the power of functional programming. However, fully understanding these topics isn't necessary when you're using functional programming.

A few years ago, when many developers started talking about functional programming (FP) as the **best way to approach concurrency**, I decided it was time to judge for myself. I expected to learn some new ideas, but I assumed I would stick with object-oriented programming (OOP) as my primary approach to software development. I was wrong.

As I learned about functional programming, I found good ideas for improving concurrency, as I expected, but I also found that it brought new **clarity** about the **design of types** and functions. It also allowed me to write more modular code. Functional programming made me **rethink where module boundaries** should be, and how to make those **modules better for reuse**. I found that the functional programming community was building innovative and more **powerful type systems** than OOP, with **correctness** as a primary goal. I also concluded that functional programming is a better fit for the unique challenges of our time, like working with massive data sets and needing to be **agile** as requirements change ever more rapidly and schedules grow ever tighter. Instead of remaining an OOP developer who tosses in some FP for sea, I decided to **write functional programs that use objects judiciously**. You could say that I stayed for the concurrency, but I stayed for the “paradigm shift.”

The funny thing is, we’ve been here before. A very similar phenomenon

then put the *same object model* in the code! Even implementation details like forms of input and output, seemed ideal for object modeling.

But let’s be clear, **both FP and OOP are tools, not panaceas**. **Each has advantages and disadvantages**. It’s easy to stick with the tried and true, even when the better way is available. Even so, it’s hard to believe that objects, which have

well in the past, could be any less valuable today, isn't it? For me, my g in functional programming isn't a repudiation of objects, which have p Rather, it's a recognition that the drawbacks of objects are harder to ignore with the programming challenges of today. The times are different than t objects were ascendant several decades ago.

Here, in brief, is why I became a functional programmer and why I beli

learn about it, too. For me, functional programming offers the best ap the following challenges, which I face every day.

I Have to Be Good at Writing Concurrent Programs

It used to be that a few of the "smart guys" on the team wrote most of code, using multithreaded concurrency, which requires carefully synch is shared, mutable state. Occasionally someone would get a midnigh

Today, even your phone has several CPU cores (or your next one will) to write robust concurrent software is no longer optional. Fortunately, programming gives you the right principles to think about concurrency and several higher-level concurrency abstractions that make the job far easi



Multithreaded programming, requiring synchronized access to sh mutable state, is the *assembly language of concurrency*.

Most Programs Are Just Data Management Problems

I work a lot with *big data* these days, mostly using the Apache Hadoop tools, built around *MapReduce* [Hadoop]. When you are ingesting te data each day, when you need to cleanse and store that data, then do

I've come to believe that faithfully representing the domain object model be questioned. *Object-relational mapping* (ORM) and similar forms of ware add overhead for transforming relational data into objects, moving around the application, then ultimately transforming them back to relational updates. Of course, all this extra code has to be tested and maintained.

I know this practice arose in part because we love objects and we often data, or maybe we just hate working with relational databases. (I speak experience.) However, relational data, such as the result sets for queries collections that can be manipulated in a functional way. Would it be directly with that data?

I'll show you how working directly with more fundamental collection

mizes the overhead of working with object models, while still avoiding promoting reuse.

Functional Programming Is More Modular

Years ago, I had a large client that struggled to get work done with the base. Their competition was running circles around them. One day I that captured their problems in a nutshell. I walked by a five-foot partition

up protocol layers that make possible all the wonderful things that we with computers.

There are no similar standards for object-based components. Various CORBA and COM had modest success, but ultimately failed for the same that objects are at the wrong level of abstraction. Concepts like “custom new, yet we can’t find a way to stop inventing a new representation for

However, if we notice that an object is fundamentally just an aggregation, we can see a way to define better standardized abstractions at lower levels, analogous to digit circuits. These standards are the **fundamental collections**, ***map***, and ***set***, along with “primitive” types like numbers and few well-known concepts (e.g., Money in a financial application).

A further aid to modularity is the **nature of functions** in functional programming: **avoid side effects**, making them **free of dependencies on other objects**, **easier to reuse** in many contexts.

The net result is that **a functional program defines abstractions** where **useful, easier to reuse, compose, and also test**.



Any arbitrarily complex object can be decomposed into “atomic” values (like primitives) and collections containing those values and other collections.

I Have to Work Faster and Faster

Development cycles are going asymptotically to zero length. That sounds good, especially if you started professional programming when I did, when projects took months, even years. However, today there are plenty of Internet sites that update their code *several times a day* and all of us are feeling the pressure to get work done quickly, without sacrificing quality, of course.

When schedules were longer, it made more sense to model your domain *exactly* and then to *implement* that domain in code. If you made a mistake, it would take a long time to correct with a new release. Today, for most projects, understanding the domain *exactly* is less important than delivering some value quickly. Our understanding of the domain will change rapidly anyway, as we and our customers discover new requirements with each deployment. If we misunderstand some aspect of the domain, we can correct those mistakes quickly when we do frequent deployments.

If careful modeling seems less important, faithfully *implementing* the model is even more suspect today than before. While Agile Software Development has improved our quality and our ability to respond to change, we need to make sure we keep our code “minimally sufficient” for the requirements today, yet flexible enough to adapt. Functional programming helps us do just that.

Functional Programming Is a Return to Simplicity

Functional programming is a return to simplicity, a return to simplicity against accidental complexity, the kind we add ourselves by our incoherent choices, as opposed to the inherent complexity of the problem domain. For example, much of the object-oriented middleware in our applications today is redundant and wasteful, in my opinion.

I know that some of these claims are provocative. I’m not trying to convince you to abandon objects altogether or to become an FP zealot. I’m trying to give you a new toolbox and a broadened perspective, so you can make more informed

Functional programming, in its “purest” sense, is rooted in how functions and values actually work in mathematics, which is different from how they work in most programming languages.

Functional programming got its start before digital computers even existed. The theoretical underpinnings of computation were developed in the 1930s by mathematicians like Alonzo Church and Haskell Curry.

In the 1930s, Alonzo Church developed the Lambda Calculus, which is a model for defining and invoking functions (called applying them). Today, the behavior of most programming languages reflect this model.

Haskell Curry (for whom the Haskell language is named) helped develop Lambda Calculus Logic, which provides an alternative theoretical basis for computation. Logic examines how combinators, which are essentially functions, combine to perform a computation. One practical application of combinators is to use them as building blocks for constructing parsers. They are also useful for representing planned computation, which can be analyzed for possible bugs and opportunities.

More recently, Category Theory has been a fruitful source of ideas for functional programming, such as ways to structure computations so that side effects (like input and output), which change the state of the “world,” are cleanly separated from computations with no side effects.

A lot of the literature on functional programming reflects its mathematical roots. It can be overwhelming if you don’t have a strong math background. In contrast, imperative programming seems more intuitive and approachable. Fortunately,

The first language to incorporate functional programming ideas was LISP, developed in the late 1950s and is the second-oldest high-level programming language after Fortran. The ML family of programming languages started in the 1970s with ML, Caml, OCaml (a hybrid object-functional language), and Microsoft's F#. F# is the best known functional language that comes closest to functional "purity" today, which was started in the early 1990s. Other recent functional languages include Haskell and Scala, both of which run on the JVM but are being ported to the .NET CLR. Today, many other languages are incorporating ideas from functional programming.

The Basic Principles of Functional Programming

languages considered *functional* languages? Functional languages share some common principles.

Avoiding Mutable State

The first principle is the use of immutable values. You might recall the *Pythagorean equation* from school, which relates the lengths of the sides of a right triangle:

$$x^2 + y^2 = z^2$$

If I give you *values* for the *variables* x and y , say $x=3$ and $y=4$, you can compute z (5 in this case). The key idea here is that values are never modified. It's a bit crazy to say $3++$, but you could start over by assigning new values to the variables.

Most programming languages don't make a clear distinction between a value (the contents of memory) and a variable that refers to it. In Java, we'll use final to indicate variable reassignment, so we get objects that are *immutable values*.

Why should we avoid mutating values? First, allowing mutable values makes multithreaded programming so difficult. If multiple threads can modify a shared value, you have to synchronize access to that value. This is quite error-prone programming that even the experts find challenging. [Goetz] says: make a value immutable, the synchronization problem disappears. Consequently, immutability is harmless, so multithreaded programming becomes far easier.

A second benefit of immutable values relates to program correctness. It's harder to understand and exhaustively test code with mutable values because mutations aren't localized to one place. Some of the most difficult bugs in distributed systems occur when state is modified non-locally, by client code that is far from where in the program.

Consider the following example, where a mutable `List` is used to hold orders:

```
public class Customer {  
    // No setter method  
    private final List<Order> orders;  
    public List<Order> getOrders() { return orders; }  
    public Customer(...) {...}  
}
```

It's reasonable that clients of `Customer` will want to view the list of `Order` objects. Unfortunately, by exposing the list through the getter method, `getOrders`, we've handed control over them! A client could modify the list without our knowledge. We could add a setter for `orders` and it is declared `final`, but these protections only protect against replacing a new list to `orders`. The list itself can still be modified.

We could work around this problem by having `getOrders` return a copy of the list instead of the list itself, by adding special accessor methods to `Customer` that provide controlled access to the `orders`. However, copying the list is expensive, especially for large lists.

What happens when the list of orders is supposed to change, but it has to stay immutable? Should we relent and make it mutable to avoid the overhead of making copies? Fortunately, we have an efficient way to copy large data structures; we'll use `Arrays.copyOf()` that aren't changing! When we add a new order to our list of orders, we only copy the rest of the list. We'll explore how in [Chapter 3](#).

Some mutability is unavoidable. All programs have to do IO. Otherwise, the program is just nothing but heat up the CPU, as a joke goes. However, functional programming encourages us to think strategically about when and where mutability is necessary. By encapsulating mutations in well-defined areas and keeping the rest of the code immutable, we improve the robustness and modularity of our code.

We still need to handle mutations in a thread-safe way. [Software Transactions](#) and the [Actor Model](#) give us this safety. We'll explore both in [Chapter 4](#).



Make your objects immutable. Declare fields `final`. Only provide getters for fields and then only when necessary. Be careful that mutable objects can still be modified. Use mutable collections carefully. "Minimize Mutability" in [\[Bloch2008\]](#) for more tips.

Functions as First-Class Values

In Java, we are accustomed to passing objects and primitive values to methods, turning them from methods, and assigning them to variables. This means that objects and primitives are *first-class values* in Java. Note that *classes themselves* are *class values*, although the reflection API offers information about class values.

Functions are not first-class values in Java. Let's clarify the difference between a *method* and a *function*.



A *method* is a block of code attached to a particular class. It can be called in the context of the class, if it's defined to be *static*, or in the context of an instance of the class. A *function* is more general. It is not attached to any particular class or object. Therefore, all instance methods are *functions* where one of the arguments is the object.

Java only has methods and methods aren't *first-class* in Java. You can't pass a method as an argument to another method, return a method from a method, or assign a method as a value to a variable.

However, most *anonymous inner classes* are effectively function "wrappers". Many methods take an instance of an interface that declares one method. Here's an example:

to the left of the "arrow" (\rightarrow) and the body of the function is to the right of the arrow. Notice how much boilerplate code this syntax removes!



The term *lambda* is another term for *anonymous function*. It comes from the use of the Greek lambda symbol λ to represent functions in *lambda calculus*.

For completeness, here is another example function type, one that takes two arguments of types $A1$ and $A2$ respectively and returns a non-void value of type R .

Closures

A *closure* is formed when the body of a function refers to one or more variables that aren't passed in as arguments or defined locally, but are defined in the enclosing scope where the function is defined. The runtime has to "close over" these variables so they are available when the function is actually executed, which can happen long after the original variables have gone out of scope! Java has limited support for closures in inner classes; they can only refer to final variables in the enclosing scope.

Higher-Order Functions

There is a special term for functions that take other functions as arguments or return them as results: *higher-order functions*. Java methods are limited to primitive types as arguments and return values, but we can mimic this feature with interfaces.

Higher-order functions are a powerful tool for building abstractions and reusable behavior. In Chapter 3, we'll show how higher-order functions allow for the customization of standard library types, like `Lists` and `Maps`, and also for the creation of new types. In fact, the *combinators* we mentioned at the beginning of this chapter are higher-order functions.

Side-Effect-Free Functions

Being able to replace a function call for a particular set of parameters with its return value is called *referential transparency*. It has a fundamental implication: a function with no side effects; the function and the corresponding return values are *interchangeable*, as far as the computation is concerned. You can represent the function with a value. Conversely, you can represent any value with a function call!

Side-effect-free functions make excellent building blocks for reuse, since they *depend on the context in which they run*. Compared to functions with side effects, they are also easier to design, comprehend, optimize, and test. Hence, they are preferred. They don't have bugs.

which encapsulates the required loop counting. We'll see other iterations in the next chapter, when we discuss operations on functional collections.

The classic functional alternative to an iterative loop is to use *recursion*. A function that passes through the function operates on the next item in the collection until a base point is reached. Recursion is also a natural fit for certain algorithms, such as traversing a tree where each branch is itself a tree.

Consider the following example, where a unit test defines a simple tree structure. Each node has a value at each node, and left and right subtrees. The `Tree` type definition is as follows:

Lazy vs. Eager Evaluation

Mathematics defines some *infinite* sets, such as the *natural numbers* (a sequence of natural numbers). They are represented symbolically. Any particular finite subset of the set is *evaluated only on demand*. We call this *lazy evaluation*. *Eager evaluation* evaluates all the values to represent all of the infinite values, which is clearly impossible.

Some languages are lazy by default, while others provide lazy data structures.


```

private static List<Integer> take(int count) {
    return doTake(emptyList(), count);
}

private static List<Integer> doTake(List<Integer> accumulator, int count) {
    if (count == ZERO)
        return accumulator;

    else
        return doTake(list(next(count - 1), accumulator), count - 1);
}

```

We start with a definition of zero, then use `next` to compute each natural

its predecessor. The `take(n)` method is a pragmatic tool for extracting a finite list of the integers. It returns a `List` of the integers from 1 to `n`. (The `List` type is discussed in [Chapter 3](#). It isn't `java.util.List`.) Note that the helper method `doTake` is tail-call recursive.

We have replaced values, integers in this case, with functions that compute values on demand, an example of the *referential transparency* we discussed earlier. A representation of infinite data structures wouldn't be possible without this property. Referential transparency and lazy evaluation require side-effect-free functions and immutable values.

Finally, lazy evaluation is useful for deferring expensive operations until they are needed, never executing them at all.

Declarative vs. Imperative Programming

Finally, functional programming is *declarative*, like mathematics, where relationships are defined. The runtime figures out how to compute finite values. The definition of the factorial function provides an example:

```

fact(0) = 1
fact(n) = n * fact(n - 1)

```

The definition relates the value of `factorial(n)` to `factorial(n-1)`, a recursive relation. The special case of `factorial(1)` terminates the recursion.

Object-oriented programming is primarily *imperative*, where we *tell* the computer specific steps to do.

To better understand the differences, consider this example, which provides both a declarative and an imperative implementation of the factorial function:

The `factorial` method implements a calculation of factorials, but its structure is more declarative than the `imperativeFactorial` method. I formatted the method to look similar to the definition of `factorial`.

The `imperativeFactorial` method uses mutable values, the loop counter, and a `result` that accumulates the calculated value. The method explicitly implements a particular algorithm. Unlike the declarative version, this method has lots of steps, making it harder to understand and keep bug free.

Declarative programming is made easier by *lazy evaluation*, because at runtime the opportunity to “understand” all the properties and relationships between values to determine the optimal way to compute values on demand. Like *lazy evaluation*, *functional programming* is largely incompatible with mutability and functions with side effects.

Designing Types

Whether you prefer *static* or *dynamic* typing, functional programming has many lessons to teach us about good type design. First, all functional languages emphasize the use of *core container types*, like *lists*, *maps*, *trees*, and *sets* for capturing and organizing data, which we’ll explore in [Chapter 3](#). Here, I want to discuss the benefits of functional thinking about types, enforcing valid values for variables, and applying rigor to type design.

of Lisp (as its name suggests). Don't confuse the following classic definition of the built-in `List` type.

As you read this code, keep a few things in mind. First, `List` is an Algebraic Data Type with structural similarities to `Option<T>`. In both cases, a common interface *protocol* of the type, and there are two concrete subtypes, one that represents "empty" and one that represents "non-empty."

Second, despite the similarities of structure, we'll introduce a few more idioms that get us closer to the requirements of a true algebraic data type.

prevent this as they can only create lists terminated by `EMPTY`. Hence, this is good behavior.



Pure functional programming uses recursion instead of loops, so a loop counter would have to be mutable.

We used a few idioms to enforce the algebraic data type constraint that the `archv` must be closed, with only two concrete types to represent all lists.

Maps

Let's talk briefly about `maps`, which associate keys with values, as in the following example:

```
Map<String,String> languageToType = new HashMap<String,String>();
languageToType.put("Java", "Object Oriented");
```

```
languageToType.put("Ruby", "Object Oriented");
languageToType.put("Clojure", "Functional");
languageToType.put("Scala", "Hybrid Object-Functional");
```

Maps don't make good *algebraic data types*, because the value of defining a "non-empty" type (or similar decomposition) is less useful. In part, due to the fact that the "obvious" implementation of `List` is strongly implied by a *tail design*.

There is no such obvious implementation of `Map`. In fact, we need flexible, different implementations for different performance goals. Instead, `Map`

Combinator Functions: The Collection Power Tools

You already think of lists, maps, etc. as “collections,” all with a set of common operations. Most collections support Java `Iterators`, too. In functional programming, three core operations that are the basis of almost all work you do with collections.

Filter

Create a `new collection`, keeping only the elements for which a `filter` returns `true`. The size of the new collection will be `less than or equal to` the original collection.

Map

Create a `new collection` where each element from the original collection is `transformed into a new value`. Both the original collection and the new collection have the `same size`. (Not to be confused with the `Map` data structure.)

Fold

Starting with a “seed” value, traverse through the collection and `use each element to build up a new final value` where each element from the original collection “contributes” to the final value. An example is summing a list of integers.

Many other common operations can be built on top of these three. Together, they form the basis for implementing concise and *composable* behaviors. Let’s see how.

```

package datastructures2;
...
public class ListModule {
    public static interface List<T> {
        ...
        public List<T> filter (Function1<T,Boolean> f);
        public <T2> List<T2> map (Function1<T,T2> f);
        public <T2> T2 foldLeft (T2 seed, Function2<T2,T,T2> f);
        public <T2> T2 foldRight (T2 seed, Function2<T,T2,T2> f);
        public void foreach (Function1Void<T> f);
    }

    public static final class NonEmptyList<T> implements List<T> {
        ...
        public List<T> filter (Function1<T,Boolean> f) {
            if (f.apply(head())) {
                return list(head(), tail().filter(f));
            } else {
                return tail().filter(f);
            }
        }

        public <T2> List<T2> map (Function1<T,T2> f) {
            return list(f.apply(head()), tail().map(f));
        }

        public <T2> T2 foldLeft (T2 seed, Function2<T2,T,T2> f) {
            return tail().foldLeft(f.apply(seed, head()), f);
        }

        public <T2> T2 foldRight (T2 seed, Function2<T,T2,T2> f) {
            return f.apply(head(), tail().foldRight(seed, f));
        }

        public void foreach (Function1Void<T> f) {
            f.apply(head());
            tail().foreach(f);
        }
    }
}

```

There are five new methods declared in the `List` interface. We need `fold`, `foldLeft` and `foldRight`, for reasons we'll discuss in a moment. And we need a `foreach` method for convenience.

Each implementation for the five new methods in `NonEmptyList` is recursive. There are no checks for the end of the recursion! The corresponding implementation in `EMPTY` terminates the recursion. This means we have eliminated the need for null tests, replacing them with object-oriented polymorphism!

Recall that the `filter` method will return a new `List`. It takes a `Function` `f` and applies `f` to each element. In `Empty`, `filter` just returns `EMPTY`. In `NonEmptyList`, the result of applying `f` to `head` (`f.apply(head())`) is `true`, then `filter` constructs a new list with `head` and the result of calling `filter` on the tail. Otherwise, `filter` discards `head` and the result of applying `filter` to the tail, thereby discarding `head`. So, `filter` terminates when it is called on an empty list.

The `map` method is slightly simpler, since it never discards an element. It uses recursion to traverse the list, applying `f` to each element and building up a new list of the results. Note that `f` is now of type `Function1<T, T2>`, because the goal is to transform the original elements of type `T` to be transformed into instances of the new type. At the same time, `EMPTY`'s `map` method calls `emptyList`, because it must return an `List<T2>`, instead of an object of the original type.

The `foldLeft` and the `foldRight` methods are the hardest to understand. They are actually the most important, as all other methods could be implemented in terms of them. We'll start with a general discussion of how these methods work, then we'll look at the implementation details.

The reason there are two versions is because they traverse the collection in different orders. In some cases, the ordering doesn't matter and the results will be different. There are other important differences we'll see later.

In a nutshell, `foldLeft` groups the elements from left to right, while `foldRight` groups them from right to left. It might help to start with an illustration of how these methods work. Suppose I have a list of the integers 1 through 4. I want to calculate the sum using `fold`. Consider the following example:

```
val list = List(1, 2, 3, 4)
```


So, we can see that `foldRight` can be used with infinite data structures, `n` elements will be evaluated.

However, `foldRight` has a drawback; it is not *tail recursive*. Why? No addition *after* the recursive call returns. The recursive call isn't the *tail of the algorithm*. The *tail-call optimization* can't be applied to `foldRight`.

However, `foldLeft` is tail recursive. Let's write the left-recursive version in the next example:

happens *before* the call, to construct the argument passed to the next call to `foldLeft`. Hence the recursion is a tail call, the last calculation done.

However, `foldLeft` can't be used for infinite data structures. There is no way to replace a call to `next` with the seed, as for `foldRight`. So, `foldLeft` will evaluate the expression, blowing up on an infinite data structure.

Now let's return to the implementations, starting with `foldLeft`. First, `foldLeft` is of type `Function2<T2, T, T2>`. The *first* `T2` type parameter represents the type of the value that we are building up a new value that could be just about anything (e.g., `String`, an `Integer` (for sums), etc. So, we have to pass a *starter or "seed"* value. A conventional name for this argument is *accumulator*, since it will *accumulate* the "result" of the work done up to a given point.

The second type parameter `T` for `f` is the type of the elements in the collection. The last type parameter `T2` is the final return type of the call to `foldLeft`. Note

call to `tail().foldRight` *after* the latter has returned. As we discussed, why `foldRight` is not tail recursive.



Consider these concise and precise definitions: `foldLeft` “is the fundamental list iterator” and `foldRight` “is the fundamental list recursive operator” [Shivers].

To end our discussion of `fold`, note that there is a similar operation called `foldLeft` is like `fold`, but the initial value of the collection is used as the seed. Hence,

Combinator Functions: The Collection

general, because the type of the result doesn’t have to be the same as collection elements. Also, unlike `fold`, `reduce` will fail if used on an empty collection since there is no “first” value!

Finally, we have `foreach`, the simplest of all these methods. Technically, `foreach` is disallowed in “pure” functional programming, because it performs operations on the collection.

Computations from simpler pieces. Combinators are arguably *the most reusable* we have in programming.



The `filter`, `map`, and `fold` functions are *combinators*, composable building blocks that let us construct complex computations from simpler pieces. They are *highly reusable*. The combination of `map` and `reduce` was the inspiration for the *MapReduce* approach to processing massive data sets [Hadoop].

Finally, recall that I implemented these functions using recursion, but `fold` and `foldLeft` avoid recursion, as in our `FunctionCombinatorTest` example. The

a value. While this may be fine for small objects, it will be too expensive for objects like long lists and large maps.

Fortunately, we can have both immutability and acceptable performance by allocating memory for what is actually changing and we share the unchanged parts of the original object. This approach is called *structure sharing*. Tree data structures provide the balance of performance and space efficiency required to do this. The abstraction might still be a `List`, a `Map`, or another data structure. The trees themselves are for the internal storage. Note that the trees themselves and their nodes are mutable. Otherwise, structure sharing will be dangerous, as mutation of one object will be seen by others that share the same substructure!

To simplify the discussion, let's use unbalanced binary trees. They provide $O(\log_2(N))$ search times (unless the tree is *really* unbalanced). Real

Software Transactional Memory

Chances are you've worked on an application with a database backend. One of the features of most relational databases is support for *ACID transactions*, an acronym for *atomicity, consistency, isolation, and durability*. The goal of ACID transactions is to prevent logical inconsistencies in a given set of related records, for example when simultaneous updates leave the set of records in an inconsistent state, or updates that are based on stale data, which could effectively erase more recent updates.

Software Transactional Memory (STM) brings transactions to locations in memory that are referenced by variables [STM] (see also [PeytonJones2007]). STM provides *durability*, because memory isn't durable (e.g., if the power is lost), but it provides the ACI, *atomicity, consistency, and isolation* in ACID.

The model in STM is to separate references to values from the values themselves. We saw this principle at work in Akka actors. In STM, a program has a reference to a value of interest. The STM framework provides a protocol for changing the value. The reference "points."

However, values themselves are *not* changed. They remain immutable. References change to point to new values. We saw in "Persistence Structures" on page 36 how the appropriate choice of implementation can provide an efficient way to make a new value from a large object without copying the parts that aren't changing. Rather, those parts are shared between the old and new