

Zigurd MEDNIEKS
Series Editor



ANDROID™
DEEP
DIVE

G. Blake MEIKE



ANDROID™
CONCURRENCY

FREE SAMPLE CHAPTER
SHARE WITH OTHERS



Android™ Concurrency

About the Android Deep Dive Series

Zigurd Mednieks, Series Editor

The Android Deep Dive Series is for intermediate and expert developers who use Android Studio and Java, but do not have comprehensive knowledge of Android system-level programming or deep knowledge of Android APIs. Readers of this series want to bolster their knowledge of fundamentally important topics.

Each book in the series stands alone and provides expertise, idioms, frameworks, and engineering approaches. They provide in-depth information, correct patterns and idioms, and ways of avoiding bugs and other problems. The books also take advantage of new Android releases, and avoid deprecated parts of the APIs.

About the Series Editor

Zigurd Mednieks is a consultant to leading OEMs, enterprises, and entrepreneurial ventures creating Android-based systems and software. Previously he was chief architect at D2 Technologies, a voice-over-IP (VoIP) technology provider, and a founder of OpenMSobile, an Android-compatibility technology company. At D2 he led engineering and product definition work for products that blended communication and social media in purpose-built embedded systems and on the Android platform. He is lead author of *Programming Android* and *Enterprise Android*.

Android™ Concurrency

G. Blake Meike

◆ Addison-Wesley

Boston • Columbus • Indianapolis • New York • San Francisco • Amsterdam • Cape Town
Dubai • London • Madrid • Milan • Munich • Paris • Montreal • Toronto • Delhi
Mexico City • São Paulo • Sidney • Hong Kong • Seoul • Singapore • Taipei • Tokyo

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The author and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions, custom cover designs, and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the U.S., please contact intlcs@pearson.com.

Visit us on the Web: informit.com/aw

Library of Congress Control Number: 2016937763

Copyright © 2016 Pearson Education, Inc.

All rights reserved. Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, request forms and the appropriate contacts within the Pearson Education Global Rights & Permissions Department, please visit www.pearsoned.com/permissions/.

Android is a trademark of Google, Inc.

ISBN-13: 978-0-13-417743-4

ISBN-10: 0-13-417743-6

Text printed in the United States on recycled paper at RR Donnelley in Crawfordsville, Indiana
First printing: June 2016

Editor-in-Chief
Greg Wiegand

Executive Editor
Laura Lewin

Development Editor
Sheri Replin

Managing Editor
Sandra Schroeder

Project Editor
Lori Lyons

Project Manager
Ellora Sengupta

Copy Editor
Abigail Manheim
Bass

Indexer
Cheryl Lenser

Proofreader
Natarajan

Technical Reviewers
Joe Bowbeer
Thomas Kountis

Editorial Assistant
Olivia Basegio

Cover Designer
Chuti Prasertsith

Compositor
codeMantra

**Graphics
Conversion**
Vived Graphics



This book is for my mother, Sally Ann Obermeyer Meike, who encouraged me to be a snoot, long before David Foster Wallace re-purposed the term



This page intentionally left blank

Contents at a Glance

Preface xiii

About the Author xvi

Acknowledgements xvii

1	Understanding Concurrency	1
2	Java Concurrency	9
3	The Android Application Model	29
4	Async Tasks and Loaders	41
5	Looper/Handler	69
6	Services, Processes, and Binder IPC	91
7	Periodic Tasks	127
8	Concurrency Tools	167
	Bibliography	183
	Index	185

NOTE: To register this product and gain access to bonus content, go to www.informit.com/register to sign in and enter the ISBN. After you register the product, a link to the additional content will be listed on your Account page, under Registered Products.

This page intentionally left blank

Table of Contents

Preface	xiii
About the Author	xvi
Acknowledgements	xvii

1 Understanding Concurrency 1

Concurrency Made Hard	1
Concurrency in Software	2
Concurrency in Hardware	3
Concurrency Made Simple	4
Threads	4
Atomic Execution	5
Visibility	6
Summary	7

2 Java Concurrency 9

Java Threads	9
The Thread Class	10
Runnables	11
Synchronization	11
Mutexes	12
Common Synchronization Errors	15
Volatile	17
Wait and Notify	20
Wait	20
Notify	22
The Concurrency Package	23
Safe Publication	23
Executors	25
Futures	26
Summary	27

3 The Android Application Model 29

Lifecycles and Components	29
Process Priority	31
Component Lifecycles	33
Android Applications as Web Apps	35

- The Android Process 36
 - Application Startup 36
 - The Android Main Thread 38
- Summary 40

4 Async Tasks and Loaders 41

- Async Task Architecture 41
 - Async Task Basics 45
 - AsyncTask Execution 48
 - AsyncTask Completion 50
- Using AsyncTasks 52
 - AsyncTask: Considered Dangerous 52
 - Getting It Right 58
- Loaders and Cursor Loaders 59
- AsyncTasks: What Went Wrong 67
- Summary 68

5 Looper/Handler 69

- Introducing the Looper/Handler 69
 - Basic Looper/Handler 71
 - Delegating Execution 73
- Some Gory Details 81
 - Handlers and Messages 81
 - Starting a Looper 84
 - The Native Looper 86
 - Scheduling and the Sync-Barrier 87
- Summary 88

6 Services, Processes, and Binder IPC 91

- Service Basics 92
 - Started Service Essentials 93
 - Bound Service Essentials 95
- Intents 96
- The Intent Service 99
- Bound Services 103
 - A Simple Bound Service 103
 - Binding A Service 105

Unbinding A Service	106
Binding Multiple Services	107
Service Lifecycle	109
Priorities and Flags	111
A Local Bound Service	112
Inter-Process Communication	114
Parcelables	116
Messengers	116
Using AIDL	120
Creating Processes	123
Binder, Considered Briefly	125
Binder Threads	125
Binder Data Transfer Limits	126
Binding to Death	126
Summary	126
7 Periodic Tasks	127
Task Characteristics	127
Thread Safety	128
Lifecycle-Aware	128
Smart Use of Process Priority	128
Power-Thrifty	129
The Scorecard	130
Timer and TimerTask	130
Looper/Handler	130
Custom Service-based Scheduler	133
Alarm Manager and Intent Service	134
The Alarm Manager Service	135
Schedulable Tasks	139
Sync-Adapters	142
Defining a Sync-adapter	143
How Sync-adapters Work	145
Implementing the Sync-adapter	147
Scoring the Sync-adapter	152

- The JobScheduler 155
 - Scheduling a Task 156
 - Running a Task 159
 - Task Implementation 161
 - Scoring the JobScheduler 165
- Summary 166

8 Concurrency Tools 167

- Static Analysis 167
 - Android Studio 168
 - Findbugs 169
- Annotations 177
 - JCIR Annotations 177
 - Support Library Annotations 178
- Assertions 179
- Conclusion and Best Practices 180

Bibliography 183

Index 185

Preface

I have had the opportunity during my years in working in this industry, to see concurrency in many contexts. When I was in school, it was a topic for dissertations. As a journeyman developer, I saw a lot of concurrent code, nearly all of it Java, in distributed back-end systems. Recently, I've had the opportunity to experience first-hand the turn to languages like Erlang and Scala, in the hope of making concurrent code easier to design and write.

I recall early in my career, being coached by a very supportive interviewer into reinventing double-checked locking. I certainly recall the furor, about a year later, when it was discovered that the double-checked locking idiom was not safe, and soon after, the first sighting of incorrect byte-code in the wild. Perhaps most surprising of all, though, I recently removed an implementation of the double-checked locking idiom from a piece of code written in 2015.

The constant, over this time, is the mystery and debate surrounding the topic. Perfectly competent novices suddenly balk or make naïve mistakes when concurrent code is necessary. Developers who are otherwise entirely reasonable sometimes disagree—occasionally quite vehemently—over the correctness of a particular piece of concurrent code. Their arguments, which may go on for hours, inevitably end up hinging on minutiae so fragile that the argument's actual winner is irrelevant.

I readily admit to feeling some kind of glee when I walked out of that early interview after having been led to double-checked locking. It was a shibboleth. I'd been initiated! Recently, I think I've seen the same kind of excitement in the faces of lecturers and their audiences as the lecturer passes on a secret: some fast and loose—and frequently downright incorrect—concurrency trick.

A really clever algorithm and a good shot of glee are wonderful things. We would all write better code, though, if we could strip some of the mystery and magic from concurrent programming. It would be great if instead of being the realm of the wizard, the correctness of a piece of concurrent code were something on which the opinions of two developers—even two developers with wildly different interests—rapidly converged.

Who Should Read This Book

This book is intended for developers with some experience with Android development.

If you are a novice developer, you will probably find some of the terms and concepts here unfamiliar. If you are a developer working on his or her very first Android app, you will probably be more concerned with simply getting familiar with the Android framework.

There are some really good books for you, already in existence, if you fit into either of these two categories. If that is the case, I encourage you to set this book aside for a while, enjoy the thrill of the steep part of a learning curve, and to come back when you have one complete Android app under your belt.

While it may sound obvious, I intend this book for reading. It is neither a cookbook nor a reference manual. I absolutely encourage you to try out the sample code. The examples are just sandboxes for experiments. Extend them. Try new experiments with them. Build your own personal understanding of the details of the Android OS. I hope, though, that you don't have to prop the book open next to your laptop to get value from it.

I was at one time bewitched by Perl. Looking back, I think that a major reason for that is that it was so much fun to read the Camel book (not the pink one, the blue one: the 2nd edition). I recall many enjoyable hours far from a keyboard and with no specific application in mind, simply reading that book.

I do not compare myself to Larry Wall. I do not by any means compare Android to Perl. I do hope, though, that you enjoy just reading this book. I hope that it is something that will make good company, perhaps on a plane flight or a long commute.

How This Book Is Organized

The first three chapters of this book will be a review for the audience for whom the book is intended.

I urge you not to skip Chapter 1, at least. It presents a model of concurrency that is somewhat atypical and that is the basis for the discussion in the following chapters.

Chapters 2 and 3 are intentionally short. They are a refresher for some basic ideas and provide an opportunity to reintroduce some common idioms. Experienced developers may choose to skim them.

Chapter 4 is a cautionary tale.

The heart of the book is Chapters 5 through 7. These chapters are a deep dive into some of the details of the Android operating system.

Chapter 8 is dessert: a bit of a how-to for some concurrency tools.

Example Code

Most of the code shown in examples in this book can be found on GitHub at <https://github.com/AndroidConcurrencyDeepDive>. If you experiment with them and discover something interesting or amusing, by all means submit a pull-request to share it with others.

Register your copy of *Android Concurrency* at informit.com for convenient access to downloads, updates, and corrections as they become available. To start the registration process, go to informit.com/register and log in or create an account. Enter the product ISBN 9780134177434 and click Submit. Once the process is complete, you will find any available bonus content under "Registered Products."

Conventions Used in This Book

The following typographical conventions are used in this book:

- **Bold** indicates new terms, URLs, email addresses, filenames, and file extensions.
- `Constant width` is used for program listings, as well as within paragraphs to refer to program elements such as variable or function names, databases, data types, environment variables, statements, and keywords.
- **Constant width bold** shows commands or other texts that should be typed by the user.
- *Constant width italic* shows texts that should be replaced with the user-supplied values or with the values determined by the context.

Note

A Note signifies a tip, suggestion, or general note.

About the Author

Blake Meike is a passionate engineer, architect, and code poet. As an author, speaker, and instructor, he has taught thousands of people how to write Android apps that aren't toys. He has more than 20 years of coding experience, most of it with Java, building systems as large as Amazon's massively scalable AutoScaling service and as small as a pre-Android OSS Linux/Java-based platform for cell phones. He is co-author of several other books, including O'Reilly's bestselling "Programming Android" and Wiley's "Enterprise Android." Blake holds a degree in Mathematics and Computer Science from Dartmouth College and was a founding member of Twitter University. He lives in Oakland, CA, and works for Cyanogen Inc.

Acknowledgements

This book owes its existence to a long list of contributors. My longtime colleague, Zigurd Mednieks, proposed it. Laura Lewin, Executive Editor at Pearson Technology Group, gave it a green light, and Carol Jelen, Literary Agent at Jelen Publishing, put a deal together.

The production staff at Pearson deftly turned the manuscript into a book. Illustrator Jenny Huckleberry cleaned up the figures and Development Editor Sheri Replin, Copy Editor Abigail Manheim Bass, and Project Editor Lori Lyons cleaned up the drafts. Editorial Assistant Olivia Basegio and Project Manager Ellora Sengupta kept us all together through the process. My thanks to all of you.

A special “thank you” and my eternal gratitude go to Laura Lewin, Executive Editor. Her patience with my complete inability to stick to any kind of schedule surely qualifies her for immediate sainthood.

On this book, my first as a solo author, my technical editors Joe Bowbeer, Thomas Kountis, and Zigurd Mednieks contributed so much that I am tempted to list them as co-authors. They were everything that I could have hoped for in a technical review team. Each, a recognized expert in the field, took the time to read carefully, understand, and then suggest changes as small as rewording and as big as re-thinking. While any remaining errors are strictly my own, any clarity and accuracy are due to them. I am incredibly lucky to have had their help. Thanks, guys!

Finally and always, thanks to my wonderful wife, Catherine, who endured yet another year of weekends with her husband on the couch, headphones on, incommunicado. That glitter heart is still for you, babe!

This page intentionally left blank

Understanding Concurrency

We propose to use the delays τ as absolute units of time which can be relied upon to synchronize the functions of various parts of the device.

John von Neumann

In order to build correct, concurrent Android programs, a developer needs a good model of concurrent processes, how they work, and what they are for. Concurrency actually isn't a big deal for most normal humans. For any multi-celled animal—arguably even for viruses—it is just normal existence. It is only those of us obsessed with computers that give a second thought to the idea of walking and chewing gum at the same time.

Concurrency Made Hard

Walking and chewing gum isn't easy in the strange world of Dr. John von Neumann. In his 1945 paper, "The First Draft Report on the EDVAC" (von Neumann 1954), he describes the architecture of one of the very first electronic digital computers. In most ways, that architecture has changed very little in the seventy years since. Throughout their history, digital computers have been, roughly speaking, gigantic balls of state that are transformed, over time, by a sequence of well-defined operations. Time and order are intrinsic parts of the definition of the machine.

Most of computer science has been the discussion of clever sequences of operations that will transform one machine state into another, more desirable, state. As modern machines commonly have more than 10^{14} possible states, those discussions are already barely manageable. If the order in which transformations take place can vary, the discussion necessarily broadens to include all possible combinations of all possible states, and becomes utterly impossible. Sequential execution is the law of the land.

Concurrency in Software

Of course, computer languages are written for humans. They are intended to help us express an algorithm (the sequence of instructions that transforms the machine state) efficiently, correctly, and, perhaps, even in a way that future human readers can understand.

Early programming languages were, essentially, an extension of the hardware. Even today many are reflections of the machine architecture they were originally designed to control. Nearly all of them are **procedural** and consist of lists of instructions for changing (**mutating**) the state of memory. Because it is simply too difficult to reason about all of the possible states of that memory, languages have, over time, become more and more restrictive about the state changes they allow a developer to express. One way to look at the history of programming language design is as a quest for a system that allows developers to express correct algorithms easily, and not express incorrect ones at all.

The very first languages were machine languages—code that translated, one-for-one, into instructions for the computer. These languages were undesirable for two reasons. First, expressing even a very simple idea might take tens of lines of code. Second, it was much too easy to express errors.

Over time, in order to restrict and manage the ways in which a program could change state, languages have narrowed the choices. Most, for instance, restrict program execution from arbitrarily skipping around between instructions to using now-familiar conditionals, loops, and procedure calls. Modules and eventually OOP (Object-Oriented Programming) followed, as ways of separating a program into small, understandable pieces and then limiting the way those pieces can interact. This modularized, building-block approach makes modern languages more abstract and expressive. Some even have well-developed type systems that help prevent errors. Almost all of them, though, are still imperative: lists of instructions for changing machine state.

Functional Programming

While most computer research and development focused on doing more and more complicated things, on bigger and faster hardware based on von Neumann architecture, a small but persistent contingent has pursued a completely different idea: *functional programming*.

A purely functional program differs from a procedural program in that it does not have mutable state. Instead of reasoning about successive changes to machine state, functional languages reason about evaluating functions at given parameters. This is a fairly radical idea and it takes some thinking to understand how it could work. If it were possible, though, it would have some very appealing aspects from a concurrency point of view. In particular, if there is no mutable state, there is no implicit time or order. If there is no implicit order, then concurrency is just an uninteresting tautology.

John McCarthy introduced Lisp, the first functional language, in 1958, only a year or two after the creation of the first commonly accepted procedural language, Fortran. Since then, Lisp and its functional relatives (Scheme, ML, Haskell, Erlang, and so on) have been variously dismissed as brilliant but impractical, as educational tools, or as shibboleths for hipster developers. Now that Moore's law (Moore, 1965) is more likely to predict the number of processors on a chip

than the speed of a single processor, people are not dismissing functional languages anymore. (By 1975, Moore formalized this concept when he revised his original thoughts, and said the number of integrated circuit (IC) components would double every two years.)

Programming in a functional style is an important strategy for concurrent programming and may become more important in the future. Java, the language of Android, does not qualify as a functional language and certainly does not support the complex type system associated with most functional languages.

Language as Contract

Functional or procedural, a programming language is an abstraction. Only a tiny fraction of developers need to get anywhere near machine language these days. Even that tiny fraction is probably writing code in a virtual instruction set, implemented by a software virtual machine, or by chip firmware. The only developers likely to understand the precise behavior of the instruction set for a particular piece of hardware, in detail, are the ones writing compilers for it.

It follows that a developer, writing a program in some particular language, expects to understand the behavior of that program by making assertions in that language. A developer reasons in the language in which the program is written—the abstraction—and almost never needs to demonstrate that a program is correct (or incorrect) by examining the actual machine code. She might reason, for instance, that something happens 14 times because the loop counter is initialized to 13, decremented each time through the loop, and the loop is terminated when the counter reaches 0.

This is important because most of our languages are imperative (not functional) abstractions. Even though hardware, registers, caches, instructions pipelines, and clock cycles typically don't come up during program design, when we reason about our programs we are, nonetheless, reasoning about sequence.

Concurrency in Hardware

It is supremely ironic that procedural languages, originally reflections of the architecture they were designed to control, no longer represent the behavior of computer hardware. Although the CPU of an early computer might have been capable of only a single operation per tick of its internal clock, all modern processors are performing multiple tasks simultaneously. It would make no sense at all to idle even a quarter of the transistors on a 4-billion gate IC while waiting for some single operation to complete.

Why Is Everything Sequential?

It is possible that sequential execution is just an inherent part of how we humans understand things. Perhaps we first imposed it on our hardware design and then perpetuated it in our language design because it is a reflection of the way our minds work.

Hardware is physical stuff. It is part of the real world, and the real world is most definitely not sequential. Modern hardware is very parallel.

In addition to running on parallel processors, modern programs are more and more frequently interacting with a wildly parallel world. The owners of even fairly ordinary feature-phones are constantly multitasking: listening to music while browsing the web, or answering the phone call that suddenly arrives. They expect the phone to keep up. At the same time, sensors, hardware buttons, touch screens, and microphones are all simultaneously sending data to programs. Maintaining the illusion of “sequentiality” is quite a feat.

A developer is in an odd position. As shown in Figure 1.1, she is building a set of instructions for a sequential abstraction that will run on a highly parallel processor for a program that will interact with a parallel world.

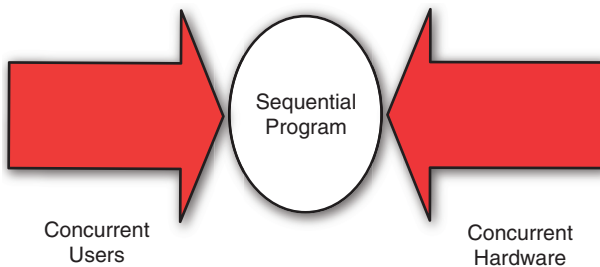


Figure 1.1 A sequential program in a concurrent world

Concurrency Made Simple

The purpose of the discussion, up to this point, has been to reframe the idea of concurrency. Concurrency is not a way to make a program run faster. It is not a complex juggling trick that ninja coders use to keep multiple balls in the air at one time. On the contrary, it is the apparent sequential execution of a program that is the complex trick. Sequential execution is an illusion maintained by a cabal of compiler writers and hardware architects. Concurrency is simply the relaxation of a fabricated constraint.

Threads

In the developer’s environment, where time and order are rigid and implicit constraints, “concurrency” is just another word for “order unspecified”—the way things are everywhere else. A concurrent program is just one that announces to an already-concurrent world that its correctness does not depend on the ordering of events that occur in separate components. In a concurrent program, those separate, partially ordered components are called **threads of execution**, or just **threads**. Within a thread, instructions are still executed in rigidly sequential order. The order of execution of instructions in two separate threads, however, is completely unspecified.

In the early days of computing, the choice of threads as a model for concurrency was not obvious. Developers that needed out-of-order processing were forced to brew their own

concurrency constructs. Both literature and code from the 1960s contain a wide variety of models for asynchronous execution.

Threads probably originated in the late 1960s, with IBM’s OS/360. They were called “tasks,” and were an OS-level service that saved developers the trouble of building their own concurrency abstraction. In 1991 Java, called Oak at the time, adopted the thread model and supported it in the language, even on operating systems that did not.

Even today, threads are not the only model for concurrency. Languages such as Erlang, Go, and Clojure, for instance, each use an entirely different model.

Introducing threads into the programming model does not present an intrinsic problem. Operating two cars in parallel causes no problems unless they both try to occupy the same space at the same time. Similarly, operating two threads that are completely independent is also perfectly safe. There are millions of programs, each concurrently running in its own thread of execution on millions of separate computers, at this very moment. Most of these programs don’t interact with each other in any way and their behavior is perfectly well defined. The problems only arise when threads need to share state and resources.

Atomic Execution

When multiple threads change state that is accessible to both, the results can easily be nondeterministic. Because there is no relationship between the order in which statements are executed, subtle changes in timing can change the result of running the program.

Consider the following code:

```
executionCount++;
someTask();
```

Just by inspection, it seems likely that the variable `executionCount` is meant to count the number of times that the function `someTask` is called. In a concurrent environment, however, this code, as it stands, does not have deterministic behavior because the `++` operation is not atomic—it is not a single, indivisible action. Table 1.1 demonstrates an execution sequence that fails to record one execution.

Table 1.1 Non-Atomic Execution

executionCount = 4	
Thread 1	Thread 2
read execution count (4)	read execution count (4)
increment (5)	increment (5)
store execution count (5)	store execution count (5)
call someTask	call someTask

Synchronization is the basic mechanism through which multiple Java threads share state in such a way that the result of the interaction is deterministic. In itself, synchronization is a simple idea: a **critical section** of code is protected by a mutual-exclusion lock or **mutex**. When a thread enters the critical section—that is, it begins executing instructions from it—it is said to seize the mutex. No other thread can enter the section until the first thread leaves.

The previous code becomes deterministic if only a single thread is allowed to execute the critical section at any given time:

```
synchronized(this) {
    executionCount++;
    someTask();
}
```

Synchronization is the crux of creating correct concurrent Java programs, and is the basis for a lot of things that are definitely not simple. Those things make up the content of the rest of this book.

Visibility

There is one more thing, however, that is simple. Remember that the previous example is an abstraction! It is written in a computer language—Java, in this case—and is, therefore, related to the actual behavior of hardware only by the grace of compiler writers, JVM developers, and hardware architects. Those two Java statements translate into hundreds of microinstructions, many of them executed in parallel, over tens of hardware clock cycles. The illusion that there are two statements, happening in order, is no more than an illusion.

Maintaining the illusion is not something that near-the-metal developers do naturally. On the contrary, they find sequential programs naive, clumsy, and wasteful. They are only too happy to fix them by re-ordering instructions, executing multiple instructions in parallel, representing a single piece of program state as multiple copies, and so on. By doing so, they do their very best to make optimal use of the immense power of the multiple processors that comprise even the tiny devices we carry in our pockets.

In general, we're glad to have them perform these optimizations. They make our programs run faster, on multiple hardware platforms, using tricks in which application developers are just not that interested. There is one important condition on this optimization, however: *They must not break the illusion of sequentiality!* In other words, compilers and hardware pipelines can reorder and parallelize all they want to optimize the code, as long as developers can't tell that they did it.

In making a program concurrent, a developer clearly states that there is no sequential dependency between the states controlled by different threads. If there is no sequential dependency, a compiler should be free to perform all sorts of optimizations that would otherwise have been unsafe. Without an explicit ordering between events in different threads, the compiler is free to make changes in the execution sequence of one thread without any consideration of the statements in any other.

A correct concurrent program is one that abides by the contract for maintaining an illusion. Negotiation between application programmers and hardware developers produce a language, and that language is a contract. The application developers get their illusion of sequential execution, something about which they can reason. The hardware developers get a toolbox of clever tricks they can use to make programs run fast. In the middle is the contract.

In Java, that contract is called the **memory model**. On one side of the contract is the application programmer, reasoning about her program in a high-level language. On the other side of the contract are the compiler writers, virtual machine developers, and hardware architects, moving everything that isn't explicitly forbidden. Developers who talk about hardware when discussing concurrency are missing the point. A correct concurrent program is not about hardware; it is about developers keeping their end of the contract.

Fortunately, in Java, the contract is easily stated. The following single sentence states it almost completely:

Whenever more than one thread accesses a given state variable, and one of them might write to it, they all must coordinate their access to it using synchronization.
(Göetz, et al. 2006)

A correct concurrent Java program is one that abides by this contract—no more, no less. Note in particular that whether a thread reads or writes mutable state does not affect its need for synchronization in any way.

Summary

Concurrency itself is nothing to be scared of. We all deal with it in the real world, all day, every day. What is difficult is writing concurrent programs for computers based on mutable state. It is difficult because the concept of order is such an important implicit basis for the way we reason about our programs.

Concurrency is the relaxation of the rigid order inherent in imperative computer languages. Java's mechanism for doing this is the thread. A thread executes instructions in an order that is not related to the order of execution of instructions in other threads. Developers use mutual exclusion locks (mutexes) to control thread access to critical sections of code, thereby limiting the number of ways that two different threads can execute instructions in the section.

Most of today's computer languages manufacture an illusion of sequential execution. Behind the scenes, however, they fiercely reorder, parallelize, and cache to make the best use of hardware. The only thing that prevents those optimizations from making a program behave non-deterministically, is a contract. A correct program is one that abides by that contract.

No one ever said that concurrency was easy. It is, however, fairly simple. Just follow the contract.

This page intentionally left blank

Index

A

abstractions, programming languages as, 3

accepting loader events (listing 4.18), 62

account type

definition (listing 7.6), 144

in sync-adapters, 143–144

Activities

in application model, 30–31

AsyncTask errors, 55–58

Services versus, 92–93

**addOnFileDescriptorEventListener
method, 86**

**administrative fields in Messages,
81–83**

**AIDL (Android Interface Definition
Language), 120–123**

**AlarmManager, scheduling tasks with
IntentService, 134–142**

AlarmManagerService, 135–139

constraints on schedulable tasks,
139–142

AlarmManagerService, 135–139

Android applications. See applications

**Android Interface Definition Language
(AIDL), 120–123**

**Android model, web apps versus desktop
apps, 67–68**

Android Studio

Findbugs, 169–175

static analysis tools, 168–169

annotations, 177

- JCIP annotations, 177–178
- support library annotations, 178–179

applications

- component lifecycles, 29–31, 33–34
- creating processes, 123–125
- definition, 30
- main thread, 36
- process priority, 31–33
- restrictions on, 29
- startup, 36–38
- web applications versus, 35

assertions, 179–180**asynchronous execution**

- AsyncTasks. *See* AsyncTasks
- with closures, 41–45
- cursor loaders, 59–66
- with loopers. *See* Looper/Handler framework

asynchronous messages, 87**AsyncTasks, 41**

- completion, 50–52
- cursor loaders and, 59–66
- doInBackground method, 45–48
- errors
 - concurrency errors, 52–55
 - lifecycle errors, 55–58
- execution, 48–49
- as Services, 99–101
- when to use
 - autonomous tasks, 59
 - cancellable tasks, 59

atomic execution, 5–6. *See also* synchronization

- single-threaded UIs, 39
- volatile keyword, 19

autonomous tasks, AsyncTasks as, 59

B

background execution. *See* asynchronous execution**best practices for concurrency, 180–181****Binders, 113, 125**

- data transfer limits, 126
- process failure, reporting, 126
- threads, 125

@BinderThread annotation, 178**binding**

- multiple Services, 107–108
- Services, 105

bindService method, 104**bound Services, 95–96, 103–104**

- binding, 105
- binding multiple Services, 107–108
- flags, 111–112
- local bound Services, 112–114
- unbinding, 106–107

bulk inserts

- initial example (listing 4.8), 51
- interruptible (listing 4.9), 51–52

C

callbacks, futures versus, 26–27**cancelable tasks, AsyncTasks as, 59****canceling AsyncTasks, 50–52****cellular radio power management, 129–130****classes**

- Messenger, 116–119
- Thread, 10–11
- Timer, 130
- TimerTask, 130

clock type in scheduling tasks, 136**Cloneable interface, 99****closures, 41–45**

clustering tasks, 129–130

Collections.unmodifiable method, 54

communication. *See* inter-process communication (IPC)

completion of **AsyncTasks**, 50–52

components

- lifecycles of, 29–31, 33–34
 - AsyncTask** errors, 55–58
- list of, 29–30

concurrency

- AsyncTask** errors, 52–55
- atomic execution, 5–6
- best practices, 180–181
- difficulty of, 1
- in hardware, 3–4
- illusion of sequential execution, 6–7
- sequential execution versus, 4
- in software
 - functional programming languages, 2–3
 - languages as contracts, 3
 - procedural programming languages, 2
- threads, 4–5. *See also* threads

Concurrency package, 23

content authority in **sync-adapters**, 143

content provider definition

- listing 7.7, 144
- listing 7.13, 152

ContentProvider, 143

contracts

- binding **Services**, 105
- maintaining sequential execution illusion, 6–7
- programming languages as, 3

copy-on-write memory sharing, 36–38

createFromParcel method, 116

cross-task communication (listing 4.7), 49

cursor loaders, 59–66

custom service-based schedulers, 133

D

data transfer limits for **Binders**, 126

database query

- with anonymous class (listing 4.4), 44
- initial example (listing 4.1), 41–42

dataset change notification

- initial example (listing 7.14), 154
- for **sync-adapters** (listing 7.15), 155

deadlock, 16–17, 38–39

declaring **Services** (listing 6.1), 92

defining **sync-adapters**, 143–145

delayed animation

- improved example (listing 5.2), 75–76
- initial example (listing 5.1), 74
- messaged example (listing 5.3), 79–80

dependency injection frameworks, 141

desktop applications versus web apps in **Android** model, 67–68

doInBackground method, 45–48, 100

E

effective immutability, 23

enqueueing **Messages** in **Looper/Handler** framework, 78

Enterprise Android (Mednieks, et al.), 143

epoll mechanism, 86

errors

- in **AsyncTasks**
 - concurrency errors, 52–55
 - lifecycle errors, 55–58
- in synchronization, 15–17

exclusion files, 177

execution

- of AsyncTasks, 48–49
- atomic execution. *See* atomic execution

executors

- futures, 26–27
- optimum number of threads, 25–26

explicit Intents, 96–97

F

fields in Messages, 81

- administrative fields, 81–83
- messaging fields, 83
- remote fields, 84

Findbugs, 169

- with Android Studio, 169–175
- with Gradle, 175–177

“The First Draft Report on the EDVAC”
(von Neumann), 1

flag field (Messages), 83

flags on bound Services, 111–112

flash memory, limitations, 29

functional programming languages, 2–3

functions, passing closures to
(listing 4.3), 43

futures

- callbacks versus, 26–27
- definition, 26

G

garbage collection, process reaping
versus, 34

Göetz, Brian, 9, 21

Gradle, 175–177

Grigorik, Ilya, 130

@GuardedBy annotation, 178

H

Handler framework. *See* Looper/Handler
framework

hardware, concurrency in, 3–4

High Performance Browser Networking
(Grigorik), 130

hints. *See* annotations

hybrid services, 93, 153

I

idealistic backgrounding (listing 4.2), 42

IllegalMonitorStateException
(listing 2.13), 21

@Immutable annotation, 178

immutable data, sharing between
threads, 23

implicit Intents, 97–98

initializing

- loaders (listing 4.17), 61
- threads as loopers, 84

inner class implementation
(listing 4.15), 57

Inspections in Android Studio, 168–169

intent filters, 97

Intents, 96–99

- in application model, 30–31
- explicit Intents, 96–97
- implicit Intents, 97–98

IntentService, 99–103

- scheduling tasks with AlarmManager,
134–142

- AlarmManagerService, 135–139
- constraints on schedulable tasks,
139–142

interfaces

- Cloneable, 99
- Parcelable, 99, 116

Runnable, 11, 73–77
 Serializable, 99
inter-process communication (IPC), 114–115
 AIDL, 120–123
 Messenger class, 116–119
 optimizing, 112–114
 Parcelable interface, 116
interruptions, 21
 bulk inserts, 51–52
 flag state, 51
 to started Services, 94–95
IPC (inter-process communication). See
inter-process communication (IPC)
isCancelled method, 50–52
isIdle method, 86

J

***Java Concurrency in Practice* (Göetz and Peierls, et al), 9, 23, 177–178**
JCIP annotations, 177–178
JobSchedule Service (listing 7.16), 156
JobScheduler, 155–165
 implementing tasks, 161–165
 running tasks, 159–161
 scoring, 165

K

keywords

 synchronized, 12, 14–15
 this, 13
 volatile, 17–20
killing processes, 33

L

LARGE_TASK_EXECUTOR, 49
lifecycle awareness in task scheduling, 128

lifecycles

 AsyncTask errors, 55–58
 of components, 29–31, 33–34
 process priority, 31–33
 of Services, 109–111

linkToDeath method, 126

listings

 accepting loader events, 62
 account type definition, 144
 AIDL definition, 120
 AIDL-generated code, 120–122
 Alarm Scheduler task execution, 140
 Android support annotations, 179
 AsyncTask Service, 100
 AsyncTasks
 concurrency errors, 52, 53, 54
 execution, 48
 initial example, 46
 lifecycle errors, 56
 local copies of mutable data structures, 54–55
 binding multiple Services to single Context, 108
 bound Service client, 103–104
 bulk inserts
 initial example, 51
 interruptible, 51–52
 content provider definition, 144, 152
 cursor loaders
 complete, 65–66
 creating, 63
 database query
 with anonymous class, 44
 initial example, 41–42
 dataset change notification
 initial example, 154
 for sync-adapters, 155

- deadlock, 17
 - declaring Services, 92
 - delayed animation
 - improved example, 75–76
 - initial example, 74
 - messed example, 79–80
 - explicit Intents, 96–97
 - Findbugs
 - failure example, 174–175
 - filter, 177
 - in Gradle, 176
 - successful example, 173
 - idealistic backgrounding, 42
 - IllegalMonitorStateException, 21
 - incorrect cross-task communication, 49
 - incorrect synchronization (single thread), 17–18
 - incorrect volatile usage, 20
 - initializing loaders, 61
 - inner class implementation, 57
 - IntentFilter, 97
 - IntentService helper method, 102–103
 - JobSchedule Service, 156
 - JobScheduler tasks
 - implementing, 161–163
 - scheduling, 157–158
 - local Service, 112–113
 - local Service client, 113
 - looper creation
 - with handler, 85
 - initial example, 84
 - without race condition, 85
 - low-power periodic scheduling with AlarmManager, 137–138
 - managed object implementation, 122–123
 - manifest for low-power scheduler, 138
 - message enqueueing, 82
 - Messenger Service, 118–119
 - Messenger usage, 117–118
 - minimal bound Service, 109
 - mutex errors, 16
 - passing closures to functions, 43
 - periodic Looper task, 131–132
 - periodic scheduling with AlarmManager, 134
 - reentrant monitors, 15
 - remote managed object usage, 123
 - safe publication, 24–25
 - Service injection, 114
 - simple account creation, 150–151
 - singleton bound Service, 111
 - skeleton Android application, 67
 - skeleton cursor list activity, 60
 - skeleton Java application, 67
 - spawning threads, 10
 - specifying component's process, 124, 125
 - sync-adapter service
 - complete service, 148
 - implementation, 148–149
 - initial example, 145
 - sync-adapters
 - definition, 143
 - implementation, 149–150
 - synchronized methods, 14
 - synchronized static methods, 14–15
 - synchronizing on objects, 12
 - synchronizing on this, 13
 - threads with Runnable, 11
 - volatile keyword, 19
 - wait and notify methods, 22
- loaders, 59–66**
- local bound Services, 112–114**

local processes, 124

locks, definition, 12. *See also* mutexes

loop method, 71

Looper/Handler framework, 69–71

- Java classes in, 71
- main thread as looper, 38–39, 74
- Messages
 - administrative fields, 81–83
 - enqueueing, 78
 - fields in, 81
 - messaging fields, 83
 - remote fields, 84
- method overloading, 78–80
- native Looper, 86
- Runnable interface, posting, 73–77
- safe publication, 69–70
- scheduling tasks, 130–133
- starting loopers, 84–86
- sync-barriers, 87–88
- task execution example, 71–73

M

main thread, 9, 36

- as looper, 38–39, 74

@MainThread annotation, 178

marshaling, 99

McCarthy, John, 2–3

McLuhan, Marshall, 167

memory leaks, 80

memory model, 7

- copy-on-write memory sharing, 36–38

MessageQueue, 70

- native Looper, 86
- sync-barriers, 87–88
- task execution example, 71–73

MessageQueue

- .addOnFileDescriptorEventListener method, 86**

MessageQueue.isIdle method, 86

MessageQueue.postSyncBarrier method, 87

MessageQueue

- .removeOnFileDescriptorEventListener method, 86**

MessageQueue.removeSyncBarrier method, 88

Messages

- enqueueing in Looper/Handler framework, 78
- fields in, 81
 - administrative fields, 81–83
 - messaging fields, 83
 - remote fields, 84
- method overloading and, 78–80

Message.setAsynchronous method, 87

messaging fields in Messages, 83

Messenger class, 116–119

methods

- addOnFileDescriptorEventListener, 86
- bindService, 104
- createFromParcel, 116
- doInBackground, 45–48, 100
- isCancelled, 50–52
- isIdle, 86
- linkToDeath, 126
- loop, 71
- notify, 22–23
- notifyAll, 22–23
- onBind, 95, 104, 105
- onCancelled, 50
- onConnected, 95
- onCreate, 74
- onHandleIntent, 101

- onPause, 76–77
- onPostExecute, 47, 50
- onProgressUpdate, 47
- onServiceConnected, 105
- onStartCommand, 93–95, 100–101
- overloading, 78–80
- postSyncBarrier, 87
- publishProgress, 47
- removeOnFileDescriptorEventListener, 86
- removeSyncBarrier, 88
- run, 76–77
- setAsynchronous, 87
- stopSelf, 101
- stopService, 101
- synchronized, 14–15
- unbindService, 104, 106–107
- unmodifiable, 54
- wait, 20–21
- writeToParcel, 116

minimal bound Service (listing 6.10), 109

model-view-controller (MVC) pattern, deadlock, 38–39

monitors. *See also* **mutexes**

- definition, 12
- reentrant monitors, 15

Moore’s law, 2–3

multiple Services, binding, 107–108

Murphy, Mark, 59

mutable data

- local copies of, 54
- sharing between threads, 23–25

mutating state, 2

mutexes, 6

- definition, 12
- errors in, 15–16
- example usage, 12–13
- reentrant monitors, 15

- synchronized methods, 14–15
- this keyword, 13
- wait method, 20–21

MVC (model-view-controller) pattern, deadlock, 38–39

N

- native Looper, 86**
- notify method, 22–23**
- notifyAll method, 22–23**
- @NotThreadSafe annotation, 178**

O

- onBind method, 95, 104, 105**
- onCancelled method, 50**
- onConnected method, 95**
- onCreate method, 74**
- onHandleIntent method, 101**
- onPause method, 76–77**
- onPostExecute method, 47, 50**
- onProgressUpdate method, 47**
- onServiceConnected method, 105**
- onStartCommand method, 93–95, 100–101**
- oom_adj attribute, 31–32**
- oom_score_adj attribute, 32**
- opportunistic suspension, 129**
- optimizing IPC (inter-process communication), 112–114**
- optimum number of threads, 25–26**
- overloading methods, 78–80**

P

- packages, Concurrency, 23**
- Parcelable interface, 99, 116**
- passing closures to functions (listing 4.3), 43**

PendingIntents, 135

periodic tasks. See **scheduling tasks**

posting Runnable interface in Looper/Handler framework, 73–77

postSyncBarrier method, 87

power usage in task scheduling, 129–130

priority of processes. See **process priority**

procedural programming languages, 2

process priority, 31–33

- flags on bound Services, 111–112
- for Services, 93
- in task scheduling, 128–129

process reaping, garbage collection versus, 34

processes

- creating, 123–125
- inter-process communication (IPC), 114–115
 - AIDL, 120–123
 - Messenger class, 116–119
 - optimizing, 112–114
 - Parcelable interface, 116
- local, 124
- remote, 124
- reporting failure, 126
- terminating, 33

programming languages

- as contracts, 3
- functional, 2–3
- illusion of sequential execution, 6–7
- procedural, 2

proxies, 115

publishProgress method, 47

R

race conditions, 13, 77

Receivers in application model, 30–31

reentrant monitors, 15

remote fields in Messages, 84

remote processes, 124

removeOnFileDescriptorEventListener method, 86

removeSyncBarrier method, 88

run method, 10, 76–77

Runnable interface, 11, 73–77

running

- JobScheduler tasks, 159–161
- sync-adapters, 145–147, 153–154

runtime concurrency checks, 179–180

S

safe publication, 23–25, 69–70

scheduling tasks

- AlarmManager and IntentService, 134–142
 - AlarmManagerService, 135–139
 - constraints on schedulable tasks, 139–142
- characteristics of, 127–128
 - lifecycle awareness, 128
 - power usage, 129–130
 - process priority, 128–129
 - scorecard for, 130
 - thread safety, 128
- custom service-based schedulers, 133
- JobScheduler, 155–165
 - implementing tasks, 161–165
 - running tasks, 159–161
 - scoring, 165
- Looper/Handler framework, 130–133
- sync-adapters, 142–155
 - defining, 143–145
 - implementing, 147–152

- running, 145–147, 153–154
 - scoring, 152–155
 - sync-barriers and, 87–88
 - Timer and TimerTask classes, 130
- scheduling window size in scheduling tasks, 136**
- scorecard (task scheduling characteristics), 130**
 - AlarmManager and IntentService, 141
 - custom service-based schedulers, 133
 - JobScheduler, 165
 - Looper/Handler framework, 132
 - sync-adapters, 155
 - Timer and TimerTask classes, 130
- sequential execution**
 - concurrency versus, 4
 - hardware versus software, 3–4
 - illusion of, 6–7
- SERIAL_EXECUTOR, 48**
- Serializable interface, 99**
- Service injection (listing 6.14), 114**
- Services**
 - Activities versus, 92–93
 - AlarmManagerService, 135–139
 - AsyncTasks as, 99–101
 - bound Services, 95–96, 103–104
 - binding, 105
 - binding multiple Services, 107–108
 - flags, 111–112
 - local bound Services, 112–114
 - unbinding, 106–107
 - custom service-based schedulers, 133
 - hybrid services, 93, 153
 - Intents, 96–99
 - explicit Intents, 96–97
 - implicit Intents, 97–98
 - IntentService, 99–103
 - scheduling tasks with AlarmManager, 134–142
 - JobSchedule Service (listing 7.16), 156
 - lifecycle of, 109–111
 - process failure, reporting, 126
 - process priority, 32
 - started Services, 93–95
 - when to use, 91
- setAsynchronous method, 87**
- sharing**
 - data between threads, 23–25
 - memory, 36–38
- single-threaded UIs, 38–39**
- singleton bound Service (listing 6.11), 111**
- singletons, 106**
- software, concurrency in**
 - functional programming languages, 2–3
 - languages as contracts, 3
 - procedural programming languages, 2
- spawning threads (listing 2.1), 10**
- start method, Thread class, 10**
- started Services, 93–95**
- starting**
 - applications, 36–38
 - loopers, 84–86
- starving threads, 19**
- static analysis tools, 167–168**
 - Android Studio, 168–169
 - Findbugs, 169
 - with Android Studio, 169–175
 - with Gradle, 175–177
- static methods, 14–15**
- stopSelf method, 101**
- stopService method, 101**
- stubs, 115**
- support library annotations, 178–179**

switched messages, 83

sync-adapters, 142–155

- defining, 143–145
- implementing, 147–152
- running, 145–147, 153–154
- scoring, 152–155

sync-barriers, 87–88

synchronization, 6, 11

- errors, 15–17
- mutexes
 - definition, 12
 - example usage, 12–13
 - reentrant monitors, 15
 - synchronized methods, 14–15
 - this keyword, 13
- notify and notifyAll methods, 22–23
- volatile keyword, 17–20
- wait method, 20–21

synchronized keyword, 12, 14–15

T

target field (Messages), 82

task clustering, 129–130

task scheduling. See scheduling tasks

terminating processes, 33

this keyword, 13

Thread class, 10–11

thread safety in task scheduling, 128

THREAD_POOL_EXECUTOR, 48–49

threads, 4–5

- asynchronous execution with closures, 41–45
- atomic execution, 5–6
- Binder threads, 125
- executors
 - futures, 26–27
 - optimum number of threads, 25–26

initializing as loopers, 84

interruptions, 51

main thread, 9, 36
 as looper, 38–39, 74

Runnable interface, 11

sharing data, safe publication, 23–25

synchronization, 11

- errors, 15–17
- mutex usage, 12–13
- notify and notifyAll methods, 22–23
- reentrant monitors, 15
- synchronized methods, 14–15
- this keyword, 13
- volatile keyword, 17–20
- wait method, 20–21

Thread class, 10–11

@ThreadSafe annotation, 178

timed-wait, 86

Timer class, 130

TimerTask class, 130

tools

- annotations, 177
 - JCIP annotations, 177–178
 - support library annotations, 178–179
- assertions, 179–180
- static analysis, 167–168
 - Android Studio, 168–169
 - Findbugs, 169–177

type-safe templates, 45

U

UI thread, 9, 36

 as looper, 38–39, 74

@UiThread annotation, 178

unbinding Services, 106–107

unbindService method, 104, 106–107
unmarshaling, 99
unmodifiable method, 54
Urdike, John, 41

V

varargs, 47
visibility, 17–20
volatile keyword, 17–20
von Neumann, John, 1

W

wait method, 20–21
wake-locks, 129, 137, 138–139
weak references, 57–58

web applications

Android applications versus, 35
in Android model, 67–68

when field (Messages), 82

Whyte, William H., 127

@WorkerThread annotation, 178

writeToParcel method, 116

Y

Yegge, Steve, 106

Z

Zygote, 36–38