

Understanding Software Dynamics

Richard L. Sites



Foreword by Luiz André Barroso, Google Fellow

FREE SAMPLE CHAPTER

SHARE WITH OTHERS



ADDISON-WESLEY PROFESSIONAL COMPUTING SERIES

Understanding Software Dynamics

The Pearson Addison-Wesley Professional Computing Series

Brian W. Kernighan, Consulting Editor



Visit informit.com/series/professionalcomputing for a complete list of available publications.

The **Pearson Addison-Wesley Professional Computing Series** was created in 1990 to provide serious programmers and networking professionals with well-written and practical reference books. Pearson Addison-Wesley is renowned for publishing accurate and authoritative books on current and cutting-edge technology, and the titles in this series will help you understand the state of the art in programming languages, operating systems, and networks.



Make sure to connect with us!
informit.com/socialconnect

Understanding Software Dynamics

Richard L. Sites

◆ Addison-Wesley

Boston • Columbus • New York • San Francisco • Amsterdam • Cape Town
Dubai • London • Madrid • Milan • Munich • Paris • Montreal • Toronto • Delhi • Mexico City
São Paulo • Sydney • 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: 2021944164

Copyright © 2022 Pearson Education, Inc.

Cover image: Art Heritage/Alamy Stock Photo

All rights reserved. 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.pearson.com/permissions.

ISBN-13: 978-0-13-758973-9

ISBN-10: 0-13-758973-5

ScoutAutomatedPrintCode

*Dedicated to the memory of Chuck Thacker,
a true Friend of the Electron who could do more
performance analysis in his head than most mortals.*

This page intentionally left blank

Contents at a Glance

Contents	ix
Foreword	xix
Preface	xxi
Acknowledgments	xxv
About the Author	xxvii

I Measurement 1

1	My Program Is Too Slow	3
2	Measuring CPUs	15
3	Measuring Memory	31
4	CPU and Memory Interaction	49
5	Measuring Disk/SSD	61
6	Measuring Networks	85
7	Disk and Network Database Interaction	111

II Observation 131

8	Logging	133
9	Aggregate Measures	141
10	Dashboards	157
11	Other Existing Tools	167
12	Traces	193
13	Observation Tool Design Principles	209

III Kernel-User Trace 217

14	KUtrace: Goals, Design, Implementation	219
15	KUtrace: Linux Kernel Patches	227
16	KUtrace: Linux Loadable Module	239
17	KUtrace: User-Mode Runtime Control	245
18	KUtrace: Postprocessing	249
19	KUtrace: Display of Software Dynamics	257

IV Reasoning 267

- 20 What to Look For 269
- 21 Executing Too Much 271
- 22 Executing Slowly 279
- 23 Waiting for CPU 289
- 24 Waiting for Memory 299
- 25 Waiting for Disk 307
- 26 Waiting for Network 319
- 27 Waiting for Locks 337
- 28 Waiting for Time 357
- 29 Waiting for Queues 361
- 30 Recap 383

A Sample Servers 387

B Trace Entries 391

Glossary 397

References 405

Index 415

Contents

Foreword	xix
Preface	xxi
Acknowledgments	xxv
About the Author	xxvii

I Measurement 1

1 My Program Is Too Slow 3

1.1 Datacenter Context	3
1.2 Datacenter Hardware	5
1.3 Datacenter Software	6
1.4 Long-Tail Latency	7
1.5 Thought Framework	9
1.6 Order-of-Magnitude Estimates	9
1.7 Why Are Transactions Slow?	11
1.8 The Five Fundamental Resources	12
1.9 Summary	12

2 Measuring CPUs 15

2.1 How We Got Here	15
2.2 Where Are We Now?	19
2.3 Measuring the Latency of an add Instruction	20
2.4 Straight-Line Code Fail	21
2.5 Simple Loop, Loop Overhead Fail, Optimizing Compiler Fail	21
2.6 Dead Variable Fail	24
2.7 Better Loop	25
2.8 Dependent Variables	26
2.9 Actual Execution Latency	26
2.10 More Nuance	27
2.11 Summary	28
Exercises	28

3 Measuring Memory 31

3.1 Memory Timing	31
3.2 About Memory	32
3.3 Cache Organization	34
3.4 Data Alignment	36
3.5 Translation Lookaside Buffer Organization	36

- 3.6 The Measurements 37
- 3.7 Measuring Cache Line Size 38
- 3.8 Problem: N+1 Prefetching 40
- 3.9 Dependent Loads 41
- 3.10 Non-random Dynamic Random-Access Memory 42
- 3.11 Measuring Total Size of Each Cache Level 43
- 3.12 Measuring Cache Associativity of Each Level 45
- 3.13 Translation Buffer Time 46
- 3.14 Cache Underutilization 46
- 3.15 Summary 46
- Exercises 47
- 4 CPU and Memory Interaction 49**
 - 4.1 Cache Interaction 49
 - 4.2 Simple Matrix Multiply Dynamics 51
 - 4.3 Estimates 51
 - 4.4 Initialization, Cross-Checking, and Observing 52
 - 4.5 Initial Results 53
 - 4.6 Faster Matrix Multiply, Transpose Method 55
 - 4.7 Faster Matrix Multiply, Subblock Method 57
 - 4.8 Cache-Aware Computation 58
 - 4.9 Summary 58
 - Exercises 59
- 5 Measuring Disk/SSD 61**
 - 5.1 About Hard Disks 62
 - 5.2 About SSDs 64
 - 5.3 Software Disk Access and On-Disk Buffering 66
 - 5.4 How Fast Is a Disk Read? 68
 - 5.5 A Little Back-of-the-Envelope Calculation 71
 - 5.6 How Fast Is a Disk Write? 72
 - 5.7 Results 73
 - 5.8 Reading from Disk 73
 - 5.9 Writing to Disk 77
 - 5.10 Reading from SSD 80
 - 5.11 Writing to SSD 82
 - 5.12 Multiple Transfers 82
 - 5.13 Summary 83
 - Exercises 84

6	Measuring Networks	85
6.1	About Ethernet	87
6.2	About Hubs, Switches, and Routers	89
6.3	About TCP/IP	89
6.4	About Packets	90
6.5	About Remote Procedure Calls (RPCs)	91
6.6	Slop	93
6.7	Observing Network Traffic	94
6.8	Sample RPC Message Definition	96
6.9	Sample Logging Design	99
6.10	Sample Client-Server System Using RPCs	100
6.11	Sample Server Program	101
6.12	Spinlocks	101
6.13	Sample Client Program	102
6.14	Measuring One Sample Client-Server RPC	105
6.15	Postprocessing RPC Logs	106
6.16	Observations	107
6.17	Summary	108
	Exercises	109
7	Disk and Network Database Interaction	111
7.1	Time Alignment	111
7.2	Multiple Clients	117
7.3	Spinlocks	118
7.4	Experiment 1	118
7.5	On-Disk Database	121
7.6	Experiment 2	121
7.7	Experiment 3	125
7.8	Logging	127
7.9	Understanding Transaction Latency Variation	128
7.10	Summary	128
	Exercises	129
II	Observation	131
8	Logging	133
8.1	Observation Tools	133
8.2	Logging	133
8.3	Basic Logging	134

- 8.4 Extended Logging 135
- 8.5 Timestamps 135
- 8.6 RPC IDs 136
- 8.7 Log File Formats 137
- 8.8 Managing Log Files 138
- 8.9 Summary 139

9 Aggregate Measures 141

- 9.1 Uniform vs. Bursty Event Rates 142
- 9.2 Measurement Intervals 143
- 9.3 Timelines 143
- 9.4 Further Summarizing of Timelines 145
- 9.5 Histogram Time Scales 147
- 9.6 Aggregating Per-Event Measurements 150
- 9.7 Patterns of Values Over Time 151
- 9.8 Update Intervals 152
- 9.9 Example Transactions 154
- 9.10 Conclusion 155

10 Dashboards 157

- 10.1 Sample Service 157
- 10.2 Sample Dashboards 159
- 10.3 Master Dashboard 159
- 10.4 Per-Instance Dashboards 163
- 10.5 Per-Server Dashboards 164
- 10.6 Sanity Checks 164
- 10.7 Summary 165
- Exercises 165

11 Other Existing Tools 167

- 11.1 Kinds of Observation Tools 167
- 11.2 Data to Observe 169
- 11.3 top Command 170
- 11.4 /proc and /sys Pseudofiles 171
- 11.5 time Command 171
- 11.6 perf Command 171
- 11.7 oprofile, CPU Profiler 173
- 11.8 strace, System Calls 176
- 11.9 ltrace, CPU C Library Calls 179

11.10	ftrace, CPU Trace	180
11.11	mtrace, Memory Malloc/Free	183
11.12	blktrace, Disk Trace	184
11.13	tcpdump and Wireshark, Network Trace	187
11.14	locktrace, Critical Section Locks	189
11.15	Offered Load, Outbound Calls, and Transaction Latency	189
11.16	Summary	191
	Exercises	191
12	Traces	193
12.1	Tracing Advantages	193
12.2	Tracing Disadvantages	194
12.3	The Three Starting Questions	194
12.4	Example: Early Program Counter Trace	197
12.5	Example: Per-Function Counts and Time	199
12.6	Case Study: Per-Function Trace of Gmail	203
12.7	Summary	207
13	Observation Tool Design Principles	209
13.1	What to Observe	209
13.2	How Frequently and For How Long?	210
13.3	How Much Overhead?	211
13.4	Design Consequences	212
13.5	Case Study: Histogram Buckets	212
13.6	Designing Data Display	214
13.7	Summary	215
III	Kernel-User Trace	217
14	KUtrace: Goals, Design, Implementation	219
14.1	Overview	219
14.2	Goals	220
14.3	Design	221
14.4	Implementation	223
14.5	Kernel Patches and Module	224
14.6	Control Program	224
14.7	Postprocessing	225
14.8	A Note on Security	225
14.9	Summary	225

15	KUtrace: Linux Kernel Patches	227
15.1	Trace Buffer Data Structures	228
15.2	Raw Traceblock Format	229
15.3	Trace Entries	230
15.4	IPC Trace Entries	232
15.5	Timestamps	233
15.6	Event Numbers	233
15.7	Nested Trace Entries	233
15.8	Code	234
15.9	Packet Tracing	234
15.10	AMD/Intel x86-64 Patches	236
15.11	Summary	237
	Exercises	237
16	KUtrace: Linux Loadable Module	239
16.1	Kernel Interface Data Structures	239
16.2	Module Load/Unload	240
16.3	Initializing and Controlling Tracing	241
16.4	Implementing Trace Calls	241
16.5	Insert1	241
16.6	InsertN	243
16.7	Switching to a New Traceblock	244
16.8	Summary	244
17	KUtrace: User-Mode Runtime Control	245
17.1	Controlling Tracing	245
17.2	Standalone ktrace_control Program	246
17.3	The Underlying ktrace_lib Library	246
17.4	The Control Interface to the Loadable Module	247
17.5	Summary	247
18	KUtrace: Postprocessing	249
18.1	Postprocessing Details	249
18.2	The rawtoevent Program	250
18.3	The eventtospans Program	251
18.4	The spantotrim Program	253
18.5	The spantospans Program	253
18.6	The samptoname_k and samptoname_u Programs	253
18.7	The makeself Program	254

18.8	KUtrace JSON Format	254
18.9	Summary	256
19	KUtrace: Display of Software Dynamics	257
19.1	Overview	257
19.2	Region 1, Controls	258
19.3	Region 2, Y-axis	259
19.4	Region 3, Timelines	260
19.5	Region 4, IPC Legend	265
19.6	Region 5, X-axis	265
19.7	Region 6, Save/Restore	265
19.8	Secondary Controls	265
19.9	Summary	266
IV	Reasoning	267
20	What to Look For	269
20.1	Overview	269
21	Executing Too Much	271
21.1	Overview	271
21.2	The Program	271
21.3	The Mystery	272
21.4	Exploring and Reasoning	273
21.5	Mystery Understood	277
21.6	Summary	277
22	Executing Slowly	279
22.1	Overview	279
22.2	The Program	279
22.3	The Mystery	280
22.4	Floating-Point Antagonist	282
22.5	Memory Antagonist	285
22.6	Mystery Understood	286
22.7	Summary	286
23	Waiting for CPU	289
23.1	The Program	289
23.2	The Mystery	289
23.3	Exploring and Reasoning	290
23.4	Mystery 2	292

23.5	Mystery 2 Understood	293
23.6	Bonus Mystery	295
23.7	Summary	297
	Exercises	297
24	Waiting for Memory	299
24.1	The Program	299
24.2	The Mystery	300
24.3	Exploring and Reasoning	300
24.4	Mystery 2: Access to a Page Table	304
24.5	Mystery 2 Understood	304
24.6	Summary	306
	Exercises	306
25	Waiting for Disk	307
25.1	The Program	307
25.2	The Mystery	307
25.3	Exploring and Reasoning	308
25.4	Reading 40MB	310
25.5	Reading Sequential 4KB Blocks	311
25.6	Reading Random 4KB Blocks	313
25.7	Writing and Sync of 40MB on SSD	314
25.8	Reading 40MB on SSD	315
25.9	Two Programs Accessing Two Files at Once	316
25.10	Mysteries Understood	317
25.11	Summary	317
	Exercises	317
26	Waiting for Network	319
26.1	Overview	319
26.2	The Programs	320
26.3	Experiment 1	321
26.4	Experiment 1 Mystery	322
26.5	Experiment 1 Exploring and Reasoning	323
26.6	Experiment 1 What About the Time Between RPCs?	327
26.7	Experiment 2	329
26.8	Experiment 3	329
26.9	Experiment 4	330
26.10	Mysteries Understood	333

26.11	Bonus Anomaly	334
26.12	Summary	336
27	Waiting for Locks	337
27.1	Overview	337
27.2	The Program	341
27.3	Experiment 1: Long Lock Hold Times	344
27.3.1	Simple Locking	344
27.3.2	Lock Saturation	345
27.4	Mysteries in Experiment 1	345
27.5	Exploring and Reasoning in Experiment 1	346
27.5.1	Lock Capture	347
27.5.2	Lock Starvation	348
27.6	Experiment 2: Fixing Lock Capture	348
27.7	Experiment 3: Fixing Lock Contention via Multiple Locks	349
27.8	Experiment 4: Fixing Lock Contention via Less Locked Work	351
27.9	Experiment 5: Fixing Lock Contention via RCU for Dashboard	353
27.10	Summary	355
28	Waiting for Time	357
28.1	Periodic Work	357
28.2	Timeouts	358
28.3	Timeslicing	358
28.4	Inline Execution Delays	359
28.5	Summary	359
29	Waiting for Queues	361
29.1	Overview	361
29.2	Request Distribution	363
29.3	Queue Structure	364
29.4	Worker Tasks	365
29.5	Primary Task	365
29.6	Dequeue	365
29.7	Enqueue	366
29.8	Spinlock	366
29.9	The “Work” Routine	367
29.10	Simple Examples	367

29.11	What Could Possibly Go Wrong?	368
29.12	CPU Frequency	369
29.13	Complex Examples	370
29.14	Waiting for CPUs: RPC Log	370
29.15	Waiting for CPUs: KUtrace	371
29.16	PlainSpinLock Flaw	374
29.17	Root Cause	375
29.18	PlainSpinLock Fixed: Observability	376
29.19	Load Balancing	377
29.20	Queue Depth: Observability	378
29.21	Spin at the End	378
29.22	One More Flaw	379
29.23	Cross-Checking	379
29.24	Summary	380
	Exercises	380
30	Recap	383
30.1	What You Learned	383
30.2	What We Haven't Covered	385
30.3	Next Steps	385
30.4	Summary (for the Entire Book)	386
A	Sample Servers	387
A.1	Sample Server Hardware	387
A.2	Connecting the Servers	388
B	Trace Entries	391
B.1	Fixed-Length Trace Entries	391
B.2	Variable-Length Trace Entries	392
B.3	Event Numbers	393
B.3.1	Events Inserted by Kernel-Mode KUtrace Patches	394
B.3.2	Events Inserted by User-Mode Code	395
B.3.3	Events Inserted by Postprocessing Code	395
	Glossary	397
	References	405
	Index	415

Foreword

Dick Sites approaches problem-solving in a way that is shockingly rare these days: he finds it almost personally offensive to make guesses, and instead he insists on understanding a phenomenon before trying to fix it. When faced with the complexity of modern computer systems, including their hardware and software, most programmers approach performance debugging armed with a hunch about what is happening and proceed to “try this, try that” with the hope that this might yield a shortcut to a solution. Those of us who use this method are implicitly giving up on the possibility of truly grasping the complex interactions that could cause a program to underperform. The idea that something computer related is beyond understanding certainly doesn’t occur to Dick. Often it is the case that basic tools that provide telemetry on a program’s behavior are missing. In those cases Dick does the obvious thing (for Dick), which is to build them, including the visualization framework that compresses essential information about program execution into readable charts that shine a bright light into program dynamics.

When you go through Dick’s remarkable career, it becomes clear why he is confident in his ability to understand complex computing systems. He became a programmer at age 10 in 1959, and his curiosity about computing resulted in a career where he studied or worked closely with giants of our field such as Fran Allen, Fred Brooks, John Cocke, Don Knuth, and Chuck Seitz, to name just a few. His accomplishments in industry are impressively broad: from co-designing the DEC Alpha Architecture to working on Adobe’s Photoshop and speeding up Google web services such as Gmail.

When I met Dick (joining DEC in 1995), he was already a legend of our field, and I had the unique pleasure of spending time with him during his Google tenure and witnessed his problem-solving approach firsthand. Readers of this book will delight in the clarity of Dick’s writing and how performance debugging problems are described as mysteries to be solved through his knowledge of hardware/software interactions and sequences of clues unveiled by observing detailed traces of program execution. This is a book that will be immensely useful for programmers and computer designers alike, in no small part because there is no other book to compare it with. It is as unique as its author.

—Luiz André Barroso, Google Fellow

This page intentionally left blank

Preface

Understanding the performance of complex software is difficult. It is even more difficult when that software is time-constrained and mysteriously exceeds its constraints now and then. Software professionals have pictures in their heads of their software's *execution dynamics*: How the various pieces work and interact together over time and estimates of how long each piece takes. (Sometimes they even document those pictures.) But when time constraints are not met, we have few tools for understanding *why*—for finding the root cause(s) of delay and other performance anomalies. This is a textbook for software developers and advanced students who work on such software.

Software dynamics refers not just to the performance or execution time of a single program thread but to the interactions between threads, between unrelated programs, and between an operating system and user programs. Delays in complex software often are caused by these interactions—code blocking and waiting for other code to wake it up, runnable code waiting for the scheduler to assign it a CPU to run on, code running slowly due to shared-hardware interference from other code, code not running at all because an interrupt routine is using its CPU, code invisibly spending much of its time in operating-system services or in page-fault handling, code waiting for I/O devices or network messages from other computers, and so on.

Time-constrained software handles repeated tasks that have periodic deadlines or tasks that have an aperiodic arrival rate of new requests each with a deadline. These tasks can have *hard* deadlines for sending control signals to moving machinery (airplanes, cars, industrial robots), *soft* deadlines such as for converting speech to text on the fly, or just *aspirational* deadlines such as for customer database lookup or web-search response times. Time-constrained also applies to phone/tablet/desktop/game user-interface responses. The term *time-constrained* is broader than the term *real-time*, which often implies hard constraints.

In each case, software tasks have a stimulus or request and a result or response. The elapsed time between the stimulus and result, the *latency* or *response time*, has some deadline. Tasks that exceed their deadlines fail, sometimes in catastrophic ways and sometimes in merely frustrating ways. You will learn how to find the root causes for these failures.

The individual tasks within such software can be called *transactions*, *queries*, *control-responses*, or *game-reactions* depending on the context. Here we will use the term *transactions* to encompass all of these. Often an end-to-end task is composed of several sub-tasks, some of which run in parallel and some of which depend on the completion of other subtasks. Sub-tasks may be CPU-bound, memory-bound, disk-bound, or network-bound. They may be executing but more slowly than expected due to interference across shared hardware resources or due to power-saving strategies in modern CPU chips. They may be waiting (i.e., not executing) for software locks or for responses from other tasks or other computers or external devices. There may be unexpected delays or interference from the underlying operating system or its kernel-mode device drivers, rather than in the programmer's user-mode code.

In many situations the software involved consists of a dozen or more *layers* or subsystems, all of which may contribute to unexpected delays and all of which may be running on separate networked computers. For example, a Google web search may spread the query across 2,000 computers, each of which does a small portion of the search and then the results are passed back and prioritized. An email message arrival in the cloud may trigger subsystems for databases, network

disk storage, indexing, locking, encryption, replication, and cross-continent transmission. An automobile-driving computer may be running 50 different programs, some of which interact on every video frame coming from a half-dozen cameras, plus radar returns, changing GPS coordinates, changing 3D acceleration forces on the vehicle, and feedback about rain, visibility, tire slippage, etc. A small database system might have query optimization and disk-access subsystems using a dozen disks spread across several networked computers. A game can have subsystems for local computation, graphics processing, and networked interactions with other players.

You will learn in this book how to design in observability, logging, and timestamps for such software, how to measure CPU/memory/disk/network behavior, how to design low-overhead observation tools, and how to reason about the resulting performance data. Once you have an accurate picture of the actual elapsed-time tasks and sub-tasks for normal transactions and also for slow ones, you can see how that reality differs from the picture in your head. At that point, substantially improving the slow transactions may take only 20 minutes of software changes. But without a good picture of reality, programmers are reduced to guessing and “trying things” to reduce long delays and improve performance. This book is about not guessing, but knowing.

All of the examples, programming exercises, and supplied software in this book are written in C or C++, based on the Linux operating system running on 64-bit AMD, ARM, or Intel processors. The reader is assumed to be familiar with developing software in this environment. We assume further that the reader has some software that is time-constrained and has performance issues that the reader wants to fix. The software should already be functional and deemed debugged, with acceptable average performance—the problem is just unexplained performance variance. The reader is assumed to have an in-the-head picture of how the software runs and can on request sketch how the pieces are *supposed* to interact in a typical transaction. Finally, the reader is assumed to know a little about CPUs, virtual memory, disk and network I/O, software locks, multi-core execution, and parallel processing. Together, we will take it from there.

We explore three major themes: Measure, Observe, and Reason.

Measure. The starting place for any performance study is to measure what is happening. A numerical measurement—transactions per second, 99th percentile response time, or dropped video frame count—tells you only what is happening but not why.

Observe. To understand why some measurement is unexpectedly slow or otherwise bad but measuring the same work again is fast, it is necessary to observe in close detail where all the time is going or what processing is being done for both normal and slow instances. For the hard case of unexpectedly bad behavior that only occurs under heavy live load, it is necessary to observe over a substantial enough time interval to have a high probability of observing several slow instances and to do so in situ with minimal distortion while running full live loads.

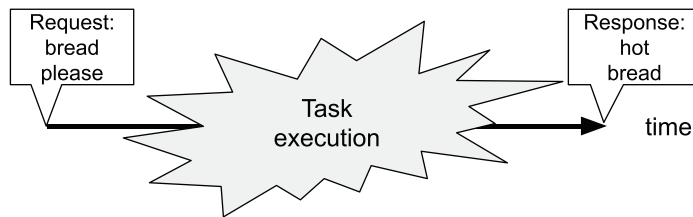
Reason (and fix). Once careful observations are available, you need to reason about what you see—how are slow instances different from normal ones, how do software and hardware interactions produce slow instances, and how can you improve the situation? In the last part of the book, we go through case-study examples of such reasoning and some of the fixes.

Following these themes, the book material is organized into four parts, including a part about building the low-overhead KUtrace observation tool:

- Part I (Chapters 1–7), **Measurement**—how to do careful measurements of the four fundamental computer resources: CPU, memory, disk/SSD, and network.

- Part II (Chapters 8–13), **Observation**—normal observation tools: logging, dashboards, counting/profiling/sampling, and tracing.
- Part III (Chapters 14–19), **Kernel-User Trace**—the design and construction of a running low-overhead Linux tracing tool that records what every CPU core is doing every nano-second, along with postprocessing programs to create dynamic HTML pages that display the resulting timelines and interactions.
- Part IV (Chapters 20–30), **Reasoning**—case studies of reasoning about the interference underlying unusual delays observed in: excess execution, slow instruction execution, waiting for CPU, memory, disk, network, software locks, queues, and timers.

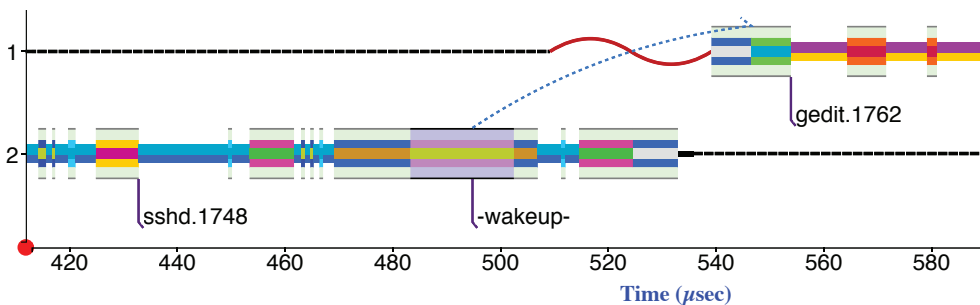
Using these ideas, you will be able to turn this picture of unexplained delay:



into the following detailed picture showing which subtasks happened when, which happened in parallel, which depended on another step finishing, and thus exactly *why* it took three hours:



The same ideas can turn an example software delay into this picture of the remote-login ssh daemon on CPU 2 waking up gedit on CPU 1:



(In Part III you will learn how to create this last kind of picture for your arbitrary software.)

This book is intended especially for engaged readers who do the included programming assignments and who implement portions of the software observation tools described.

Layered throughout this book are comments about modern complex processor chips and their performance-enhancing mechanisms. Accidentally defeating these mechanisms can create surprising delays. The careful reader will gain a deeper understanding of computer architecture and microarchitecture, along with everything else.

This is a textbook for software professionals and advanced students. But it also covers material of interest to computer hardware architects, operating system developers, system-architecture IT professionals, real-time system designers, and game developers. Its focus on understanding user-facing latency will develop skills that enhance any programmer's career.

Accessing the Source Code

The book uses several computer programs: `mystery1`, `mystery2`, and so forth. The source code for these programs is available for download from Addison-Wesley at informit.com/title/9780137589739.

Register your copy of *Understanding Software Dynamics* on the InformIT site for convenient access to updates and/or 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 (9780137589739) and click Submit. Look on the Registered Products tab for an Access Bonus Content link next to this product, and follow that link to access any available bonus materials. If you would like to be notified of exclusive offers on new editions and updates, please check the box to receive email from us.

Acknowledgments

Many people have helped along the path to this book. Amer Diwan, V. Bruce Hunt, Richard Kaufmann, and Hal Murry have actively read and provided feedback on the text. Connor Sites-Bowen, J. Craig Mudge, Jim Maurer, and Rik Farrow provided thoughtful reviews and encouragement for earlier versions and related articles. Brian Kernighan did a thorough reading of the manuscript and made suggestions to materially improve the final product.

Much of the material here was developed from graduate courses I taught after retiring from Google in 2016. I am grateful for the opportunities and the student feedback arranged by Michael Brown at the National University of Singapore; Jim Larus and Willy Zwaenepoel at Ecole Polytechnique Federale de Lausanne; Christos Kozyrakis at Stanford University; and Kevin Jeffay and Fred Brooks at the University of North Carolina.

Joshua Bakita, Drew Gallatin, and Hal Murray have done ports of KUtrace to different Unix flavors. Jim Keller and Pete Bannon provided me the opportunity to do a port at Tesla Motors. Sandhya Dwarkadas asked the key question about detecting cache interference that led to my adding instructions-per-cycle counting to KUtrace.

My early career became focused on CPU performance and tracing through the influence and guidance of Elaine Bond, Pat Goldberg, Ray Hedberg, Fran Allen, and John Cocke at IBM; Don Knuth at Stanford; and Joel Emer, Anita Borg, and Sharon Perl at Digital Equipment Corporation.

My wife of 37 years, Lucey Bowen, has been especially gracious and supportive while I spent too much time focused on completing the book.

My editor, Greg Doench, has been particularly helpful in bringing this project to a smooth completion. He took time in the early months to arrange trial runs of importing text and the extensive figures into the publishing workflow, saving time and grief near the end of the process. My copy editor, Kim Wimpsett, did a fantastic job inserting literally thousands of small improvements.

—Richard L. Sites, September 2021

This page intentionally left blank

About the Author

Richard L. Sites wrote his first computer program in 1959 and has spent most of his career at the boundary between hardware and software, with a particular interest in CPU/software performance interactions. His past work includes VAX microcode, DEC Alpha co-architect, and inventing the performance counters found in nearly all processors today. He has done low-overhead microcode and software tracing at DEC, Adobe, Google, and Tesla. Dr. Sites earned his PhD at Stanford in 1974; he holds 66 patents and is a member of the US National Academy of Engineering.

This page intentionally left blank

Chapter 6

Measuring Networks

The fourth fundamental shared resource to measure is network activity—how long do real network transmissions take and what are their dynamics? In contrast to measurements of the internal dynamics of single operations in the previous chapters, we will look at multiple overlapping network requests. The environment for disk measurements (and for that matter CPU and memory measurements) is fairly simple, as shown in Figure 6.1. There is just the one program running on one CPU and accessing one disk with a single transfer at a time.

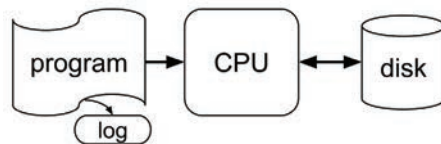


Figure 6.1 Environment for disk measurements

But the environment for network measurements is substantially more complicated, as shown in Figure 6.2. There are multiple client programs sending request messages to several server programs, which in turn send responses. These are all running on several different computers with network connections in between. Common server programs include database software.

In general, the different computers of Figure 6.2 could be located anywhere in the world, but we will concentrate in this chapter on computers that are physically close to each other, such as all inside a single datacenter room. The network connections could be Ethernet, Infiniband, Fibre Channel, or other choices, but we will concentrate on Ethernet connections. Various network protocols could be used, such as virtual channels, User Datagram Protocol (UDP), or Transmission Control Protocol/Internet Protocol (TCP/IP) software. We will concentrate on TCP/IP links, a common choice within datacenters.

The request messages and their responses could be structured in various ways; we will concentrate on *remote procedure call* (RPC) messages. Each RPC request message specifies a server computer to perform some work, the particular method (i.e., function or procedure name) to be called, and copies of all the method arguments. Each response message specifies the client computer to receive the response and the response data itself. The request and response messages can

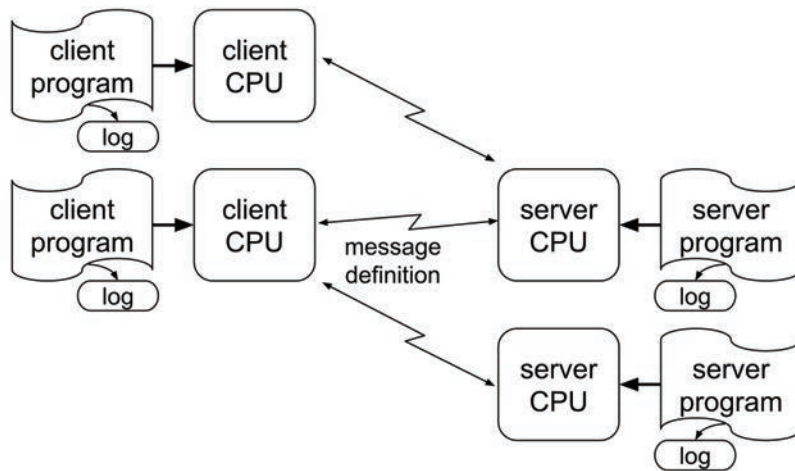


Figure 6.2 More complicated environment for network measurements

vary widely in size, from about 100 bytes to tens of megabytes. RPCs are usually *asynchronous*, meaning that the caller need not wait for the RPC response, but can instead continue executing and issuing other RPC requests in parallel, eventually waiting for responses that come back in arbitrary order. It is the highly parallel execution of many small pieces of work that allows datacenter software to respond quickly. Unlike TCP and other network protocols, RPC message formats are not standardized. This book uses a simple made-up format, described in Section 6.8.

In a large datacenter of 20,000 computers, with each computer running many different programs, an individual computer may have 10,000 network connections open at once, exchanging RPCs over all of them. While Figure 6.2 shows multiple point-to-point connections between client CPUs and server CPUs, these are just conceptual. The physical network may have just a single Ethernet link between each computer and a network router, with all the RPC traffic shared across these links. These underlying physical links and their associated kernel software are the shared network resource that we measure in this chapter.

A note on notation. The word *server* is somewhat overloaded in the computer industry. It can refer to a box of hardware that is a computer, or it can refer to a program that performs some specific function on behalf of various client programs. To add to the confusion, a server program performing a specific function is often called a *service*. In this book, when the context is not clear we will refer to *server CPU* or *sample server* for a box of hardware and *server program* for a piece of software providing some service. The unqualified term *server* will generally mean a CPU.

As discussed in Chapter 1, datacenter software consists of layers and layers of subsystems, many running in parallel and often on hundreds or thousands of different servers. All this activity is tied together with some form of network message passing or RPCs. In this chapter we will observe and measure some simple RPCs, and then in the next chapter we will measure multiple

overlapping RPCs. There are several layers of software involved, including user code, the operating system, and the TCP (transmission control protocol) stack on the client computer and the same three on the server computer. We will use RPCs from one user-mode program to another and back, measuring the behavior and delays between sample servers.

6.1 About Ethernet

Ethernet is the standard networking technology worldwide and is heavily used in datacenters. The original Ethernet at Xerox PARC in 1973 used a single coaxial cable (one wire inside a tube of a second wire with insulation in between), so it was a shared medium. Individual Alto computers connected to the wire with a *vampire tap* that poked an insulated spike through the outer wire to touch the inner wire, plus a second connection to the outer wire, as shown in Figure 6.3. (The vampire taps were shown to be unreliable and were soon superseded.) Just as polite people do when talking in a group, a computer desiring to transmit would listen to the coax (carrier-sense) waiting until it was idle and then try to transmit. During transmission, it continued to listen to determine whether the bits it transmitted were on the wire or whether they were garbled because some other computer was also starting to transmit. When that happened, both would stop transmitting, each wait a random amount of time, and then try again. Any node connected to the shared coax can observe all the packets, not just those addressed to that node. This is useful for monitoring network performance and debugging network problems, but it raises security issues.

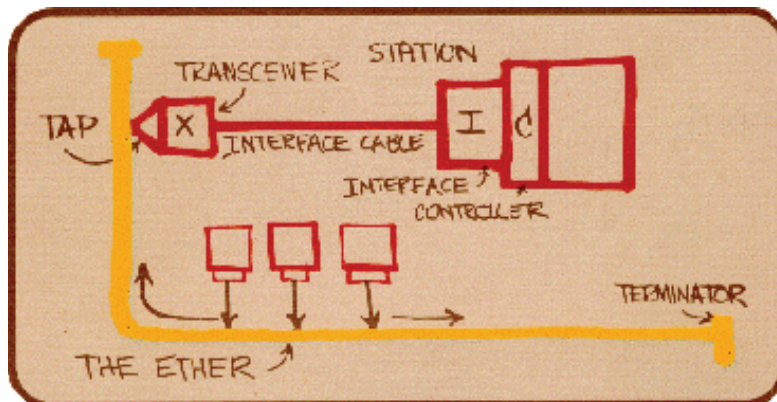


Figure 6.3 Metcalfe's original Ethernet diagram, photographed by Boggs [Metcalfe 1976]

Today, Ethernet data is transmitted as packets of up to 1,518 bytes (jumbo packets can be bigger) with gaps [Wikipedia 2021n] in between and a checksum at the end, as in Figure 6.4. Network software turns longer messages into sequences of packets. Individual packets are delivered with high probability, but are not 100% guaranteed to arrive. In particular, switches and routers that are overloaded are free to drop packets at any time. Packets with bad checksums are also discarded.

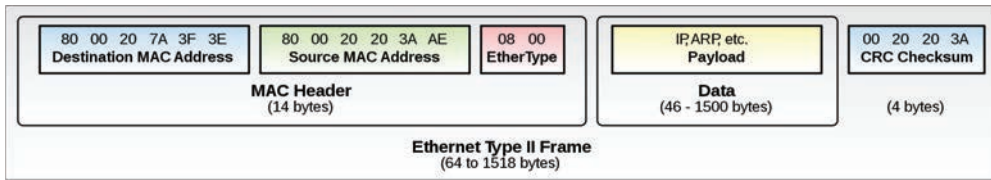


Figure 6.4 Ethernet type II frame [Wikimedia 2020a]

Each Ethernet packet starts with the 48-bit destination *Media Access Control* (MAC) address, followed by the 48-bit source MAC address, followed by a 16-bit *Ethertype* field and then the rest of the packet. The last 24 bits of a MAC address is an assigned Organizationally Unique Identifier (OUI) [IEEE 2021].

The remainder of a packet typically has several *headers* for layers of different switching protocols and then finally some user data. We will be using the TCP/IPv4 protocol pair, with a 20-byte IPv4 header giving the 4-byte IP addresses of the source and destination machines and a 20-byte TCP header giving the 2-byte port numbers on those machines plus data *sequence numbers* (SEQ) and *acknowledgment bits* (ACK) for accomplishing in-order guaranteed delivery.

We are using IP version 4 (IPv4) in our examples, but all the 32-bit IP numbers in this protocol are now used up worldwide, so the newer IPv6 is also being used in datacenters. IPv6 has 128-bit IP addresses, and an entire IPv6 header is 40 bytes instead of 20.

The MAC address is a unique 48-bit identifier assigned to each network interface controller in the world. (The original 3 Mb/sec Ethernet used 8-bit addresses.) The EtherType field specifies how to interpret the following data bytes. For TCP/IP traffic, the MAC header EtherType specifies IPv4, followed by a 20-byte IPv4 header in the first few data bytes of Figure 6.4. The IP header in turn specifies that it is followed immediately by a TCP header, which specifies that it is followed immediately by some number *N* of user message data or *payload* bytes.

While the original 3 Mb/sec Ethernet used a single shared coax cable for connections, later implementations more often use twisted-pair copper wires or optical fibers running point-to-point from each computer to a hub or switch or router. These connections have progressively increased in speed from 10 Mb/sec to 100 Mb/sec, 1 Gb/sec, 10 Gb/sec, and now 100 Gb/sec with 400 Gb/sec on the horizon, *five orders of magnitude* faster than the original.

Note that network transmission rate is traditionally measured in bits per second while disk transfer rate is traditionally measured in bytes per second. Lower-case “Mb” is megabits while upper-case “MB” is megabytes. Marketing literature often confuses these, introducing factor-of-eight errors. Deliberately quoting disk transfers in Mb/sec is a cheap way to make your numbers 8x larger.

6.2 About Hubs, Switches, and Routers

Point-to-point Ethernet connections between more than two machines require some form of switching fabric. There are three different kinds that you may encounter.

A *hub* with N links is a very cheap and now rarely used design that just reproduces one incoming transmission on all inactive outgoing links. If two or more links have incoming transmissions at once, only one is copied, and the others are dropped. Since a hub can copy only one transmission at a time, it is a shared resource like the original coax.

Switches with N links store packets at each incoming port and immediately forward them to one or more outgoing ports. Smarter switches keep tables of which MAC address destinations are attached to which port and forward only to the right destination port. A switch may store as little as two packets per incoming port, forwarding one while a second arrives. If multiple packets on different incoming ports have the same outgoing port and there is not enough buffering for all of them, some of the packets are dropped. As mentioned, Ethernet does not guarantee packet delivery, just best-effort.

A *router* is a more complex form of switch, using not only the MAC address in each packet but also higher-level IP and other header address information to select the output port for each packet. Routers are often connected to other routers so that a packet may go from one end node through several routers to another end node:

A ⇒ Router1 ⇒ Router2 ⇒ Router3 ⇒ B

Typical use in a datacenter is to have 40–50 servers mounted in a vertical rack with a router at the top (or middle) each rack. Traffic between servers within that rack is delivered directly from the top-of-rack router, while traffic destined for other racks is sent from the source top-of-rack router to one of several intermediate routers that eventually send the packet to the destination top-of-rack router and on to the destination server. Often in this case, the cross-router links run at a higher speed than the individual server links: for example, 10 Gb/sec copper-wire server links within a rack and 100 Gb/sec fiber-optic cross-rack links. We will use the phrase *on the wire* to refer to bit transmission over any kind of link. Routers often have several packets of buffering per input port, so can handle a modest amount of network congestion with several input packets destined for the same output port.

Our sample servers each have a 1 Gb/sec Ethernet port, and several are connected via a five-port switch, four ports for up to four sample servers and the fifth port connected to the rest of the building, as shown in Appendix A.

6.3 About TCP/IP

The TCP/IP design allows packets to be routed not just within a single building but anywhere in the world that is connected to the global Internet. This routing sends packets across various media—not just Ethernet links, but also long-haul dedicated fibers, radio links to satellites, WiFi connections within houses, and many more kinds of sub-networks. The complex dynamics and delays of long-distance communication are beyond the scope of this book; we will concentrate just on the complex-enough dynamics and delays of Ethernet connections within a single building.

For a message from machine A to machine B, sending software on A for *guaranteed-delivery* protocols such as TCP keeps track of packets sent that do not arrive and retransmits them. Packets are therefore not guaranteed to arrive in the order originally sent, so receiving software further tracks them and reassembles messages in receive buffers. This tracking is done by the receiving TCP software on B sending back an ACK indication to A for one or more received packets. ACKs can be sent in short packets of their own but are usually piggybacked as part of other packets already going back B \Rightarrow A. Senders have a limited number of multiple *packets outstanding*—sent but not yet acknowledged. When this limit is reached, the sender must wait until some ACKs arrive. If a packet ACK does not arrive within some configured timeout period, the sender is responsible for retransmitting that packet.

We are using TCP/IPv4 to send RPC messages between servers, with each message possibly requiring many packets. Our remote procedure calls depend on the guaranteed delivery mechanism of TCP to deliver an entire message with the pieces in proper order.

It is unlikely in our little sample server cluster that we will see packets dropped and retransmitted because of hardware errors, but we will soon try to create enough network congestion to force some packet drops because of overloaded switch buffering. To try to protect the rest of any building network from also becoming overloaded when we do saturation experiments, it is best for our lab machines to connect directly with each other through their own local switch, as described earlier.

TCP establishes a reliable connection to carry a *pair* of byte streams between two programs on two machines, one stream in each direction. These are the bi-directional connections shown in Figure 6.2. Each machine is specified by its IP address and the specific program by its port number. Two-byte port numbers range from 0..65535, but ports below 1024 are restricted to specific uses. We will use ports 12345..12348 on our sample servers for RPC traffic. (Our lab machines may also have a software firewall that closes traffic on most other ports.)

Once a connection is established, there is a stream of data available in each direction between the two machines. A machine can send an (almost) arbitrary number of bytes at once, and the TCP software deals with breaking up long messages into multiple packets, or packing multiple short messages or fragments of longer messages into single packets. The communication model is just a stream of bytes, so a given RPC message may start and end in the middle of packets.

On the other end, a machine can request receiving an (almost) arbitrary number of bytes into a buffer, but the number actually delivered at once can be less than the buffer size. The receive calls normally return any data available, up to the requested amount, rather than waiting for receipt of the full amount requested. This design allows the receiving software some flexibility in managing buffers and in managing how long to wait for data (or what else to do in the meantime). The receiving logic thus must be prepared to do multiple receive calls to get all the pieces of a single complete message and must also be prepared to receive multiple messages and partial messages at each call.

6.4 About Packets

In addition to IP and TCP headers, datacenter packets may contain additional headers. For example, virtual local area networks (VLANs) can be implemented by having a 4-byte VLAN header before the IP header. Cooperating routers deliver packets based on their VLAN header, with the

effect that packets from one virtual LAN can be prevented from reaching ports associated with other virtual LANs. This design allows multiple completely unrelated networks to use shared switching equipment. Packets without a VLAN header can be dropped by a router or sent to specific unsecured ports. Incoming packets with the wrong VLAN header for a particular port can be dropped. The goal is that each type of traffic is completely unable to observe any of the other traffic, even if some connected computers are spoofing their MAC and IP addresses to try to read, and even modify and forward, others' data. If the routers themselves operate correctly, this can give some level of security and privacy.

One use of VLANs is for a building-wide network with specific authorized machines (by MAC address) attached to specific router ports and using VLAN headers. An unauthorized machine connected to the network is not allowed to use any VLAN headers, and all it can see is a tiny default network consisting of itself and a gateway/authorization computer that may choose to stop all communication with the device, may convert it into an authorized node that can use VLANs, or may allow it to connect to an outside Internet port, thus supporting devices from guests visiting the building but otherwise allowing only limited access for those unauthorized machines.

Packets may also be encrypted. Enough initial information is left unencrypted to allow the packets to be routed, and then an encapsulation header is used to signal that the remaining bytes are to be passed on unchanged and uninterpreted by any routing mechanism. The encapsulated data can be encrypted in various ways by the sender and decrypted by the receiver. The encapsulation technique can also be used to carry byte streams that actually contain non-Internet bytes and use completely different routing protocols for some private network that connects different locations via encapsulated traffic sent over the regular Internet.

We will consider only unencapsulated packets for the rest of this book, since we are focusing on server-to-server network performance and not on all the possible ways to use the Internet.

6.5 About Remote Procedure Calls (RPCs)

Our lab experiments will use a form of remote procedure call. For a *local* procedure call, routine A calls some Method with arguments and gets back a return value, with all the code running on a single machine:

```
routine A {  
    ...  
    foo = Method(arguments);  
    ...  
}
```

For a *remote* procedure call, the idea is the same, but the Method (e.g., a C function) runs on a remote computer.

The Method name and arguments are passed to the remote server in a request message, and the return value is eventually passed back in a response message, as shown in Figure 6.5. The client and server programs are constructed with calls to an RPC library. Building, sending, and parsing the request and response messages is done by the library routines, implementing a particular RPC design. Non-blocking RPCs allow multiple RPC requests to be outstanding at once and allow responses to return out of order.

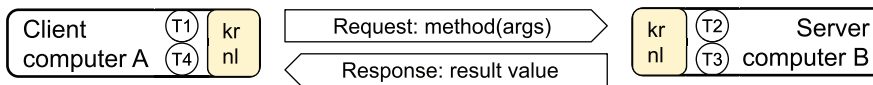


Figure 6.5 A single RPC sending a request message and eventually receiving a response message; “krnl” is kernel code; T1-T4 are timestamps in user-mode code to send/receive RPC request and response messages.

Each message is a network transmission. The request message goes from

- a user-mode client program on computer A at time T1 to
- kernel-mode code on A,
- over the network,
- to kernel-mode code on computer B,
- to a user-mode server program on B at time T2.

The response message travels in the opposite direction, at times T3 and T4. RPC latency is measured from the time T1 that the user-mode client program on A sends the request to the time T4 that the user-mode client program on A receives the response. When a response is delayed, the delay can be anywhere on the picture—request or response, user code or kernel code, machine A or machine B, send or receive network hardware. The four times T1..T4 help observe where the overall time went.

To examine the performance effects of network RPCs we will use timelines with events T1, T2, T3, and T4 indicating the RPC timing. We will draw individual RPCs as timelines with notches showing the times T1..T4, as shown in Figure 6.6. The notches do not take up much diagram space, but the human eye is quite good at picking them out, even when there are hundreds of RPC lines close together. The total RPC latency, as observed by the client user-mode program, is T4-T1. The total server time for the RPC is T3-T2.



Figure 6.6 Diagram of one RPC, showing the four times. T1 to T2 is the time from client user-mode code sending an RPC request to server user-mode code receiving that request. T2 to T3 is the server time spent performing the request. T3 to T4 is the time from server user-mode code sending the RPC response to client user-mode code receiving that response. Times T1 and T4 are taken from the client CPU’s time-of-day clock, while T2 and T3 are from the server CPU’s time-of-day clock. The two clocks may be offset from each other by microseconds to milliseconds. We will deal with clock alignment in the next chapter. w1 is the time the client kernel-mode code sends the request to the network hardware (“w” for “wire”), and w3 is the time the server kernel-mode code sends the response to the network hardware.

The return value from an RPC may be a single status number or may be thousands of bytes of data. It is convenient to always return both an overall status for the call (success, failure, specific error codes) and a possibly empty byte string of additional results.

Most datacenter software uses RPCs to send work between servers. For example, passing a paragraph of text to Google Translate via its web-page interface may send that paragraph to a load-balancing server that in turn forwards it to a least-busy translation server, which in turn may break the paragraph into sentences and send the individual sentences in parallel to a few dozen sentence servers that do sequences of multi-word phrase lookups in the source language and map into the best-score sequence out of many possible phrases in the target language. These results are then gathered back together by the translation server into a single translated paragraph.

6.6 Slop

The *slop*, or unidentified communication time, is $(T4 - T1) - (T3 - T2) = (T2 - T1) + (T4 - T3)$. When the client RPC latency and the server time are nearly equal, the slop is small. When there are communication delays (usually in the kernel code on one machine or the other, not in the network hardware), the slop can be large. Figure 6.6 shows a large slop, with the overall RPC latency about 1.5x the server time. In Chapter 7, we will also subtract the estimated transmission time for request and response messages:

$$\text{slop} = (T4 - T1) - (T3 - T2) - \text{requestTx} - \text{responseTx}$$

When the slop is large, that means that there is significant delay somewhere between the two communicating user-mode programs. In Chapter 15 we will introduce recording the $w1$ and $w3$ times of RPC headers on the wire, shown in gray in Figure 6.6. They indicate here that request and response messages hit the wire almost immediately when sent, so the long delays are in kernel code on the receiving end.

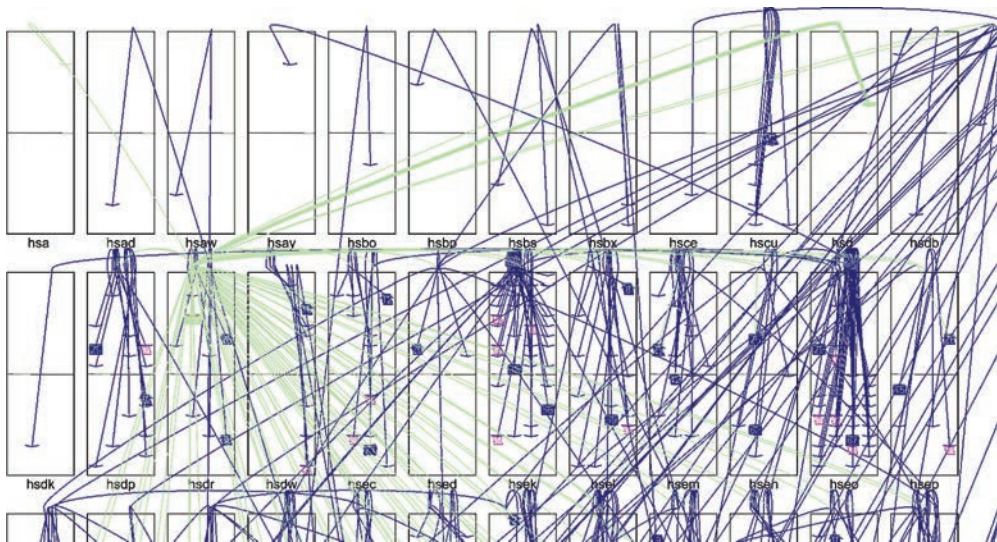


Figure 6.7 Two-level RPC call tree for a single web search, farming out one search to about 2,000 partial searches of different parts of a web index. Each rectangle represents a rack of ~50 server machines. The light green arcs show about 100 top-level RPCs from one machine in rack “hsdr” to 100 others. The dark blue arcs show about 20 second-level RPCs from each of those 100 to a total of about 2000 servers.

To complete the RPC picture, Figure 6.7 shows a two-level call tree of RPCs doing a single web search. The top-level ~100 light green RPC arcs coming from the top of the rack labeled “hsdr” are all done in parallel, and all the groups of second-level ~20 dark blue RPC arcs are also done in parallel, quickly spreading the work across about 2,000 servers. Each leaf in the call tree does a portion of the search, and those partial results are combined when all parallel RPCs in a group have returned.

6.7 Observing Network Traffic

In this chapter we will observe and measure some simple RPCs. In contrast to observing local CPU, memory, and disk activity, it takes two connected machines and two sets of software to observe network traffic. Rather than just observing isolated packets, we will observe an RPC system that has client software, server software, multi-packet RPC request and response messages, multiple server threads, and overlapped client calls. As usual, we wish to observe in enough detail to detect anomalous dynamics.

Figure 6.8 shows one example of such dynamics, captured by RPC logs and Dapper [Sigelman 2010]. Using the style of our single-RPC diagram from Figure 6.6, the notched lines in Figure 6.8 show the time layout of 93 parallel RPCs, similar to the top-level RPCs shown as timeless light green arcs in Figure 6.7.

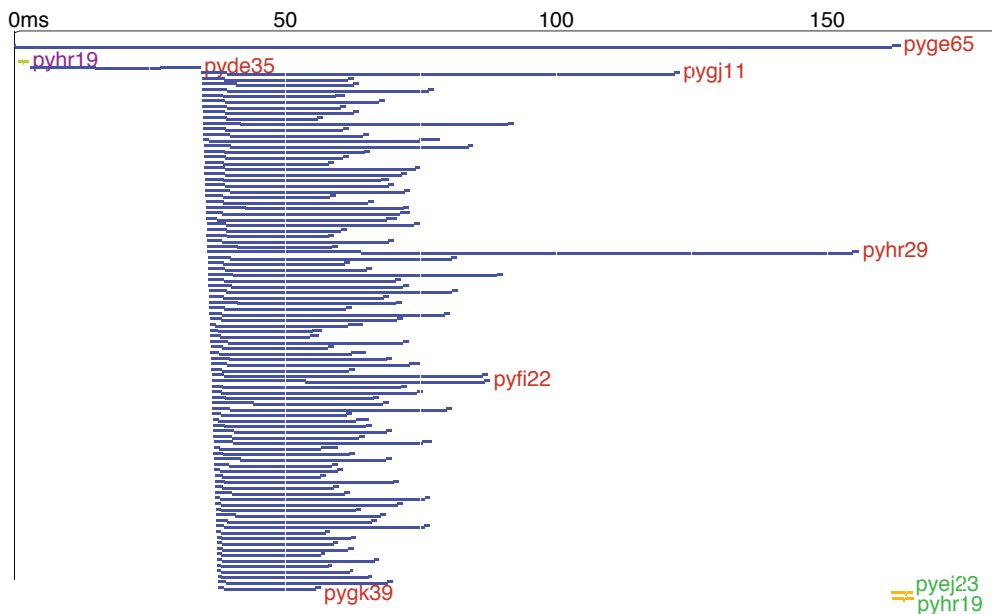


Figure 6.8 Diagram of ~100 RPCs at top level of a single web-search RPC. The “pyxxxx” notations to the right of each line are individual server names.

The very top notched blue line labeled `pyge65` in Figure 6.8 shows an incoming single RPC requesting a single web search on server `pyge65`. It takes about 160 msec total. Underneath it are the outbound sub-RPCs that this initial call spawns, directed at other servers. You can see the sub-RPC transaction latency variation and can see that the 99th percentile slowest parallel RPC on `pyhr29` determines the response time of the overall web-search RPC. You can also see that understanding and removing the sources of long latency can speed up this example by about a factor of 2, from 160 msec total to about 80 msec.

Look at the spawned RPCs. First, at the far upper left there is a short barely visible call to server `pyhr19` (yellow line just under the “m” in “0ms”) to check for a previously cached immediate answer. Using cached previous search answers noticeably speeds up identical searches. Then a blue call to `pyde35` is a *canary request* [Dean 2010]. Only when that returns successfully, i.e., without crashing `pyde35`, are the other RPCs done. If you have a request that hits a code bug that crashes a server, *and you will*, the canary strategy results in crashing just one server instead of thousands. In Figure 6.8 the canary returns, so the 90-odd parallel calls to `pygj11` .. `pygk39` are started. Not shown are the 20-odd RPCs that each of these spawn in turn, about 2,000 in total, corresponding to the dark blue arcs in Figure 6.7.

Only when `pyhr29` returns, the slowest of these parallel calls, does the initial web-search RPC complete. At the very lower right are two parallel calls to update duplicate cached results on `pyej23` and on `pyhr19` (yellow lines). These actually occur after the initial RPC completes. The vertical white line at 50 msec is just a time grid.

If you look carefully at the canary call to `pyde35`, you will notice that the request message takes over 10 msec to go from client user-mode code to server user-mode code, and the response message also takes over 10 msec to go from server to client. This slop is much longer than the slop time for most of the subsequent RPCs, so we have our first unexpected source of excess latency. For datacenter networks within a single building, the delay through the routers of the hardware switching fabric rarely exceeds 20 usec. So a delay that is 500x longer than that can only be a software, not hardware, delay, somewhere on the client or server in either user code or kernel code. We examine such delays in Part IV.

If you look carefully, the 93 parallel calls do not all start at exactly the same time—there is a slight tilt to the nearly vertical left edge. Their start times increment by about 6 usec each, reflecting the CPU time to create and send each RPC. The leftmost notch on each line shows that almost all the RPC requests arrive at their corresponding server program fairly quickly, except the calls to `pyhr29` and `pyfi22`, which take over 20 msec to arrive. This is another latency mystery to be resolved.

The rightmost notch on each line shows that almost all the RPC responses are sent soon before they arrive at the client program, so there is no latency mystery for those.

Initially, however, we would be more interested in the exceptionally slow response times of `pyhr29` and `pygj11`, since they delay the overall response time of the initial RPC by about 70 msec. Understanding those delays requires observing what is happening on each of those CPUs, applying our observation tools and thought to each in turn. The same kind of transaction log files that were used to create Figure 6.8 can be used on the logs from `pyhr29` and `pygj11` to see the dynamics of their delays. The general pattern is: observe to focus on the big issues and ignore the inconsequential artifacts, examine the important artifacts in more detail, resolve them, and then repeat.

The good news is that our picture of the RPC activity on just one machine has revealed two message-delivery latency mysteries and has pinpointed exactly the other two machines and time of day to the microsecond that contribute to overall slow response time for this one web search. Looking at multiple such web searches over a few tens of seconds will reveal whether pyhr29 and pygj11 are always slow or just happened to be slow for this one observation.

Our goal in this chapter is to capture enough information about each RPC to be able to draw diagrams like Figure 6.8 and then use those to track down root causes for delays. In later chapters, especially Chapter 26, we will add tools for observing the underlying reasons for delay(s) that our RPC diagrams reveal.

But before exploring the observation of network dynamics, we need to describe the sample “database” RPC system in a little more detail.

6.8 Sample RPC Message Definition

Local procedure calls have relatively simple dynamics—on a given CPU core, procedure A calls procedure B, and that CPU core then executes instructions in B; no further instructions in A execute until B returns. Local procedure calls may be nested, with A calling B, which in turn calls C. But these all execute sequentially on a single CPU core. We can observe a complete local call tree by capturing the entry and exit times of each procedure. Call nesting is implied by nested entry/exit times. In a multiple-core, multi-threaded-program environment, multiple local calls to B can occur simultaneously, but they are from different callers executing on different software threads possibly using different CPU cores.

Remote procedure calls are more complicated. Unlike a subroutine call, the transmission of the request message from client machine A to server machine B is not instantaneous, nor is the transmission of the response message back from B to A. Since these messages use shared network resources, other network traffic may delay them, so we at least want to capture send/receive times for each message.

As shown in Figure 6.8, an RPC may be *non-blocking*—the caller A may proceed with additional execution in parallel with B and issue additional parallel RPCs to C, D, E, etc. Eventual response messages from B, C, D, E, ... arrive at A asynchronously and not necessarily in order. To match up multiple request and response pairs in the RPC library code, each outstanding RPC is given a unique ID, and that is included in its request and response messages. To match up one RPC with any sub-RPCs that it does, each sub-RPC also includes the RPC ID of its parent; this allows us to reconstruct entire call trees.

A caller may wait for all its RPC responses before finishing, or it may finish early as in Figure 6.8 (pyge65 at the very top right returns before the calls to pyej23 and pyhr19 at the very bottom right). If a network link goes down or a server crashes, some responses may never arrive; the caller needs to detect and deal with this rather than waiting forever.

During the time that A is waiting for a response from B, other clients on the same or different machines may also be sending RPCs to B, and B may be working on those and not on A's. If that happens, a sub-RPC from B to some other server Z may be part of the work for A or for any other client of B. The parent ID shows the proper association.

If A calls B with RPCID 1234 in the request message and B subsequently calls Z on behalf of A's request, the call to Z would have parent ID 1234 and its own RPCID of perhaps 5678.

All these complications happen fairly often in a large datacenter environment.

We can observe the dynamics of a complete remote call tree only by capturing the caller/callee pairs and send/receive times for each request and response message and explicitly recording the parent caller for all nested RPCs. Most of this information must be transmitted between machines in each of the request and response messages.

For our sample RPC system, each request or response message starts with an RPC marker followed by an RPC header followed optionally by a byte string that contains the argument values for a request or the result values for a response, as shown in Figure 6.9. Each complete message is broken up into the payload data carried in one or more TCP/IP packets. We will focus in the rest of this chapter on complete messages instead of individual packets.

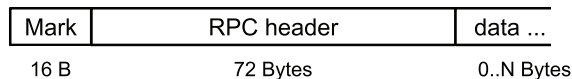


Figure 6.9 Overall structure of a request or response message in our sample RPC design

The 16-byte RPC marker, as shown in Figure 6.10, serves several purposes: delimiting messages, defining variable lengths, and sanity check.

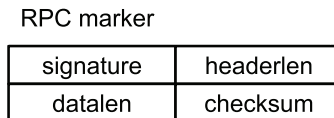


Figure 6.10 RPC marker of 16 bytes

The signature is a fixed 32-bit value. This allows a quick check that the subsequent bytes could begin an RPC message and are not something else. If somehow a TCP connection gets out of sync, it also allows scanning forward until a signature is found as a way to resynchronize. (This is not necessarily a good idea; it may be better to drop the connection and force a clean restart.) In Chapter 15, we use the signature field to filter packets that appear to be the beginning of an RPC message, recording KUtrace entries for each.

The 32-bit headerlen field gives the byte length of the following RPC header, whose size will likely vary over several months or years in a real datacenter as the RPC library is updated and expanded. To improve validity checking, headerlen values are required to be less than 2^{12} . For our sample RPC design, headerlen is always 72.

The 32-bit datalen field gives the byte length of the optional argument or result byte string that follows the RPC header. A length of 0 indicates no string. To improve validity checking and to make huge messages invalid, datalen values are required to be less than 2^{24} . The two length fields allow an RPC library to break a message into its variable-length pieces.

Finally, the 32-bit checksum field is a simple arithmetic function of the previous three fields, allowing a robust sanity check that the marker and subsequent bytes are highly likely to be the start of a valid RPC message.

The RPC marker is designed to be part of a complete network message but is not visible to the callers of the RPC library software. Those callers deal only with the RPC header and data string.

The RPC header, shown in Figure 6.11, has all the information to describe a single RPC request or response message. Fields are initialized to zero and are filled in incrementally by RPC library as an RPC is processed. For example, T1 and the first L are filled in by the RPC library when an RPC request message is about to be sent by a client program. The second L is filled in when an RPC response message is about to be sent by a server program, and T4 is not filled in until the RPC response message is received by the client program.

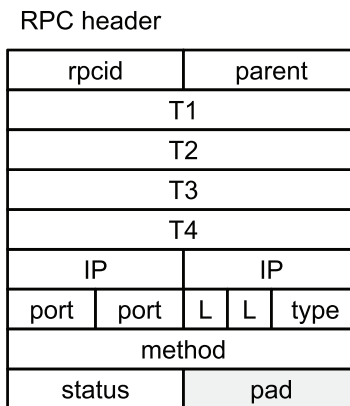


Figure 6.11 RPC header of 72 bytes.

Briefly, the naturally aligned fields are

- **RPCID** 32 bits, containing a unique ID number for each outstanding request.
- **Parent ID** 32 bits, containing the RPCID of the request that spawned the current request.
- **T1..T4** 64-bit wall-clock timestamps with microsecond resolution, giving respectively the request send time, request receive time, response send time, and response receive time; T1 and T4 are based on the client machine's time-of-day clock; T2 and T3 are based on the server machine's time-of-day clock.
- **IP** 32 bits and **port** 16 bits, giving the client and server machines' TCP/IP addresses,
- **L L** 8 bits each, giving the logarithm of the byte lengths of request and response messages.
- Message **type** 16 bits, to indicate request or response or other types of message.

- **Method** 64 bits (8 bytes), ASCII name of the routine being called, zero padded.
- **Status** 32 bits, return-value status indicating success, failure, or specific error number.
- **Padding** 32 bits, to make the header a multiple of eight bytes in total length.

The sizes of the RPC header fields are somewhat arbitrary; different sizes would work equally well. Reducing the byte lengths to logarithms is just an example of trading resolution (e.g., within 10%) for space.

This header format is somewhat less flexible than those used in real datacenters, but is sufficient for our sample RPC work.

6.9 Sample Logging Design

The client-server programs described in the coming text each write a log file of all the RPCs they process. This logging is our designed-in observability for the dynamics of the RPC system. As such, it is important that the logging is not so slow or bulky that it consumes significant resources or distorts the performance of the underlying service.

Our sample design target is to be able to process and therefore log up to 10,000 RPCs per second with little overhead. This is the right order of magnitude for a real datacenter service.

Back-of-the-envelope: If each log entry is 1,000 bytes, logging 10,000 RPC/sec would write 10MB/sec to a log file, or 864GB/day per service, with multiple services running on each server. This quickly gets bulky and can also consume significant bandwidth to a disk holding multiple log files. Each service would nearly fill up a 1TB disk every day, likely requiring multiple disk drives just for logging.

Instead, if each log entry is about 100 bytes, one service would log about 1MB/sec, and total about 86 GB/day. This is still a bit bulky, but several services writing 1MB/sec to a logging disk is an easily sustained rate, and those services would need just a single 1TB disk to hold a day's worth of log entries—a manageable amount.

For slower-rate services that handle only about 1,000 RPCs per second, we could roughly afford 1,000 bytes per log entry, although most logging does not need so much data per RPC. Do the back-of-the-envelope arithmetic at design time to document affordable limits on the logging overhead.

For our sample RPC design, the binary log format is just a copy of the current RPC header, with the full data length from the RPC mark moved to just in front of the data, and the data itself truncated or zero-extended to 24 bytes, as shown in Figure 6.12. Each log entry is thus exactly 96 bytes.

Log-system performance is a key consideration. Truncating the data keeps each log entry size bounded, and it also provides a bit of privacy for possible user-supplied data such as an email message. Including even a little of the data, though, helps identify what is going on when there is unusual latency. Recording log entries as binary fields instead of printed ASCII values saves file space and also saves CPU time for formatting all the numbers. This approach substantially reduces the logging overhead. Conversion to readable ASCII can be done by postprocessing binary log files as needed later.

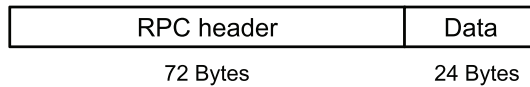


Figure 6.12 Sample log entry format, 96 bytes total

Each server program will write a local binary file of these log entries as it runs, and each client program will write its own local binary file of log entries. For an overall service running on 2,000 machines, the server programs would write 2,000 local log files on those machines. Scattered clients would write local log files on their own machines.

In a large datacenter, all programs would periodically close their log files and open new ones. A background service would gather closed log files into a single place or into a distributed file system so that they could be postprocessed efficiently. To conserve space, most logs would be thrown away after a few days. There is little point keeping data if no one will ever read it again. For our sample environment, none of that log management is done: we just deal with multiple local log files written locally on individual servers.

6.10 Sample Client-Server System Using RPCs

The supplied `server4.cc` and `client4.cc` implement a sample RPC system with logging. In addition, the supplied `dumplogfile4.cc` turns binary logs into JSON-formatted ASCII so you can see what they contain and so that they can be displayed easily.

The service provided by these programs is an in-memory key-value store (a simple database). The server program accepts RPCs that read or write key-value pairs in RAM, while the client sends such requests. The server-implemented methods are

- **Ping**(data), returns a copy of data with no key-value action
- **Write**(key, value), remembers the pair
- **Read**(key), returns value
- **Checksum**(key), returns an eight-byte checksum of the value
- **Delete**(key), removes key and its value
- **Stats**(), returns some statistics on server usage
- **Reset**(), erases all key-value pairs
- **Quit**(), stops the server and all its threads

This pair of programs should allow you to get into trouble in myriad ways. Sending a burst of 100 values of 1MB each will saturate a sample server's 1Gb/sec Ethernet for at least a second. Doing two such actions independently between different machines should overload a sample four-port switch. Sending a burst of 100,000 one-byte values will saturate CPUs quite nicely and also clog up the logging system. Overlapping a heavy-duty burst with some more modest work will likely interfere with the modest work.

6.11 Sample Server Program

The server program shown in Figure 6.13 continuously accepts RPC request messages, processes them, and sends response messages. It usually runs with multiple threads handling independent RPCs applied to a single shared database.

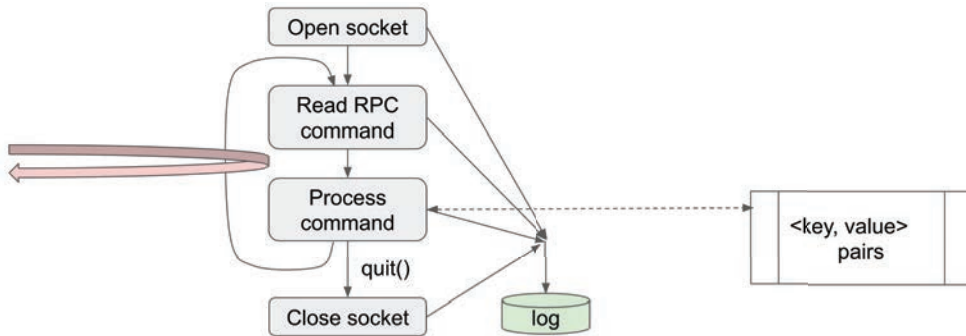


Figure 6.13 Sample server program

The program, `server4.cc`, takes two command-line arguments specifying the range of ports to listen on. For each port, it forks a dedicated listening thread. Each such thread opens a TCP/IP socket on its port and waits for RPCs, which it then executes sequentially as they arrive. The default behavior is use four sequential ports, 12345..12348, and to fork four corresponding threads.

As a safety move, each launch of `server4` will self-destruct after four minutes, to protect against runaway or zombie programs, even if `server4` is launched in the background via the command-line ampersand.

Multiple copies of `server4` can be launched, so long as they use non-overlapping port numbers. Copies of `server4` can run on multiple server machines. Nothing in the simple design except chance prevents multiple people from launching interfering runs at the same time.

Keys are byte strings restricted to less than 256 bytes, while values are byte strings restricted to less than 1.25MB (i.e., $5 * 256 * 1024$ bytes). Total RAM storage space is restricted to be less than 200MB.

The database of key-value pairs is a C++ map of strings. It is shared across all server execution threads, so each operation that touches the database takes out a simple spinlock before accessing the data. By design, this is a somewhat flawed approach—everything works, but there can be severe blocking dynamics, which we will shortly observe.

6.12 Spinlocks

The spinlock protecting a software *critical section* is our fifth shared resource, along with the four hardware resources CPU, memory, disk, and network. There are many forms of software locks, with spinlock the simplest—whenever a thread cannot acquire the lock for a critical section of

code, it simply loops (spins) trying over and over to acquire the lock, until eventually other threads free the lock and the subject thread successfully grabs it. Chapter 27 discusses locks in more detail.

The sample server code defines a C++ `SpinLock` class, which acquires a spinlock in its constructor and frees the lock in its destructor. Thus the code pattern

```
LockAndHist some_lock_name;
...
{
    SpinLock sp(some_lock_name);
    <critical section code here>
    ...
}
```

makes the inner block a critical section that can be executed by only one thread at a time.

The C++ constructor/destructor mechanism for `SpinLock` guarantees to acquire the lock `some_lock_name` upon entry to the block and to release it upon exit from the block, even for an unexpected or exception exit. This design completely removes one source of programming error—processes that sometimes fail to release a lock.

The spinlock implementation also defines a small histogram to record lock acquisition times. This is another piece of designed-in observability. A common issue with software locks is that under some circumstances a program thread has to wait much too long to acquire a lock, resulting in long transaction latency on one thread whenever *another thread* holds a lock too long. A small histogram of lock-acquire times for each lock can tell you the normal time taken to acquire a contended lock and also can show how many much-longer times occur. If there are no long acquisition times for a given lock, then lock-waiting is not a cause for a long transaction latency, and you can look elsewhere.

The locking pseudocode looks like this:

```
start = __rdtsc()
test-and-set loop to get lock
stop = __rdtsc()
elapsed_usec = (stop - start) / cyclesperusec
hist[Floorlg(elapsed_usec)]++
```

where `rdtsc` reads the x86 cycle counter and `Floorlg(x)` takes $\text{floor}(\log_2(x))$, returning 0..31 for a 32-bit unsigned int `x`. The variable `hist` is a small array of counts. Logarithm base 2 is sufficient resolution to put long and short acquisition delays into different count buckets. The `stats` command mentioned earlier returns this histogram array.

6.13 Sample Client Program

The client program, shown in Figure 6.14, takes command-line arguments and sends one or more RPCs to a specified server and port, repeated and spaced out in time in a stylized way.

Multiple instances of `client4` can be launched in the background so that they overlap in time on a single machine, and of course multiple instances can be run from several different machines. Instances of `server4` and `client4` can also run on the same machine, but will not use the network for local communication; for this degenerate case, the kernel network code moves message bytes in RAM.

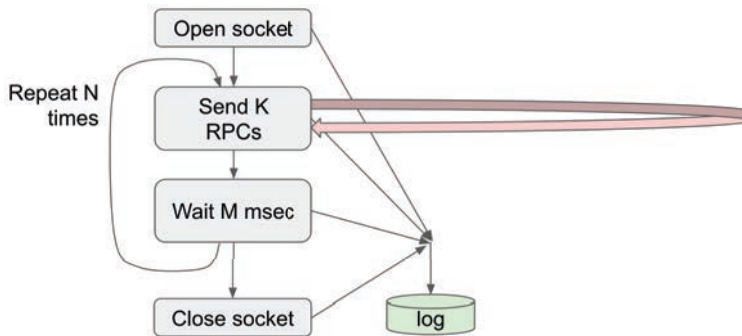


Figure 6.14 Sample client program

The command-line arguments to client4 allow specifying Rep repetitions of the pair of actions

<send K RPCs, wait M msec>

The first RPC of a burst of K is specified by its method and initial key/value data. Either field can be padded with pseudorandom data to a specified byte length. Subsequent RPCs in a burst keep the method and base strings but supply possibly incremented keys/values. Within a burst, each RPC is sent as soon as the previous one returns a response, but not before.

The client4 program takes a set of command-line parameters that form a little language:

```

./client4 server port
  [-rep number] [-k number] [-waitms number] [-seed1]
  [-verbose] command
  [-key "keybase" [+] [padlen]]
  [-value "valuebase" [+] [padlen]]
  
```

The write, read, and delete commands require a key, and the ping and write commands require a value.

-rep	Says to repeat the outer loop some number of times.
-k	Says to repeat the command some number of times (the inner loop) and then wait before continuing the outer loop.
-waitms	Says how long to wait (in milliseconds) after each burst of k commands.
-seed1	Says to seed the random number generator for padding bytes to exactly the value 1, allowing reproducible pseudorandom values.
-verbose	Prints a little about each request and response message.
-key "keybase" [+] [padlen]	Specifies a base character string that is optionally incremented and optionally padded with random characters. In the presence of -rep and -k repetitions, "+" indicates incrementing the base string, and padlen gives the padded length.
-value "valuebase" [+] [padlen]	Uses the same algorithm.

For example,

```
./client4 target_server 12345 -k 5 ping -value "vvvvv" + 10
```

sends five ping commands to target_server:port, with value strings such as

```
vvvvv_0u5j
vvvvw_trce
vvvvx_qxo1
vvvvy_1bv3
vvvvz_dglw
```

where the + specifies incrementing v w x y z in the base string, and the 10 specifies padding with random characters out to 10 characters total. Incrementing increases the low character of the base string by 1, wrapping 9 to 0, z to a, and Z to A, carrying into higher character places as needed (the next base value above would be vvvwa). Incrementing is most useful for keys and padding is most useful for values.

The individual commands are defined in a little more detail here.

ping [-value "valuebase" [+] [padlen]]

Sends an RPC request to the specified server:port containing RPC marker, RPC header, and optionally RPC data containing the specified value. The server responds with the same data.

write -key "keybase" [+] [padlen] -value "valuebase" [+] [padlen]

Sends an RPC request to the specified server:port containing RPC marker, RPC header, and RPC data containing the specified <key, value> pair. The server saves each <key, value> pair and responds with a status code, typically SUCCESS.

read -key "keybase" [+] [padlen]

Sends an RPC request to the specified server:port containing RPC marker, RPC header, and RPC data containing the specified key. The server responds with the matching value and a status code, typically SUCCESS.

delete -key "keybase" [+] [padlen]

Sends an RPC request to the specified server:port containing RPC marker, RPC header, and RPC data containing the specified key. The server deletes each matching <key, value> pair and responds with a status code, typically SUCCESS.

stats

Sends an RPC request to the specified server:port containing RPC marker, RPC header, with no RPC data. The server responds with an arbitrary status string and a status code, typically SUCCESS. For server4, the status string is a text version of the spinlock histogram described earlier—32 counts separated by single spaces.

reset

Sends an RPC request to the specified server:port containing RPC marker, RPC header, with no RPC data. The server removes all its <key, value> pairs and responds with a status code, typically SUCCESS.

quit

Sends an RPC request to the specified server:port containing RPC marker, RPC header, with no RPC data. The server responds with a status code, typically SUCCESS, and exits immediately afterward.

The client4 program prints the observed round-trip time for the first 20 RPCs it issues. At the end, it prints a log2 histogram of those times, plus total RPCs, total msec elapsed, total MB transmitted and received, and RPC messages transmitted and received per second. In a datacenter system, this is the kind of information that would be displayed on a dashboard web page.

6.14 Measuring One Sample Client-Server RPC

Recall our thought framework from Chapter 1. In thinking about performance issues, we will follow the programmer's discipline of first estimating how long some work should take, then observing how long it actually does take, and then reasoning about any differences. Figure 6.15 shows this framework again.

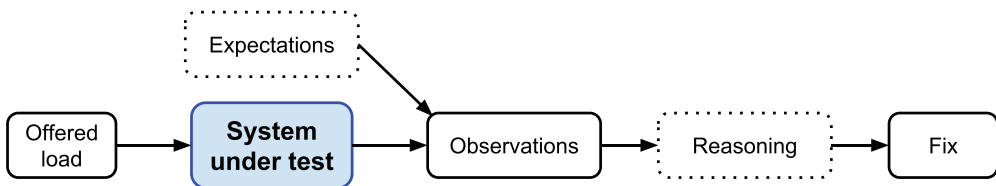


Figure 6.15 Framework for examining the performance of complex software

Consider the offered load of sending a single **write** RPC from client program to server program, with a 5-byte key and 1,000,000-byte value. What do we expect the timing to look like on our sample server configuration with a 1Gb/sec network?

Remember that 1 Gb/sec is roughly 100 MB/sec when we convert from bits to bytes and include some overhead, or about 100 bytes per microsecond. So sending ~1,000,000 bytes of RPC request should take about 10,000 usec (10 msec) across the wire. On the server side, creating a string of 1,000,000 bytes and putting it into a C++ map should take only about 100 usec or so if the main memory can transfer 10GB/sec. A short reply of 100 bytes should take about 1 usec on the wire plus some small software overhead. Overall, we might expect something like Figure 6.16 when we sketch out a timeline.



Figure 6.16 Expected RPC timing sketch for Write() of 1MB of data

This sketch is somewhat distorted, but you get the idea—a long time to transmit 1MB in the request message, 100x less time to put it into the key/value store, and 100x less time again to transmit the response message. Let's see what happens.

Compile `server4.cc` and run it with no arguments to get the default configuration of listening on four ports.

```
./server4
```

Compile `client4.cc` and run it on a different machine with these arguments:

```
./client4 target_server 12345 write -key "kkkkk" \
-value "vvvvv" 1000000
./client4 target_server 12345 quit
```

Compile `dumplogfile4.cc` and run it on the client, pointing to the 1MB write log file written by `client4`.

```
./dumplogfile4 client4_20190420_145721_dclab-1_10479.log \
"Write 1MB" \
>client4_20190420_145721_dclab-1_10479.json
```

Compile `makeself.cc` and run it on the client, pointing to the previous JSON file.

```
./makeself client4_20190420_145721_dclab-1_10479.json \
show_rpc_2019.html \
> client4_20190420_145721_dclab-1_10479.html
```

Display the resulting HTML file.

```
google-chrome client4_20190420_145721_dclab-1_10479.html
```

6.15 Postprocessing RPC Logs

The programs `client4.cc` and `server4.cc` write binary log files of 96-byte records as described earlier.

The program `dumplogfile4.cc` reads these log files and turns them into JSON files with the timestamps and other information turned into ASCII text. The JSON file has a stylized header that contains among other things the start minute of the log records, the title from the second command-line argument ("Write 1MB" previously), and some axis labels. This header is followed by lines of text for the log records. By default, only log records for receipt of a response are included, i.e., just those records that describe a full round-trip transaction. The `-all` flag includes all records.

The program `makeself.cc` reads a JSON file and writes an HTML file based on a template, incorporating the JSON information. This is the same `makeself.cc` program we encountered in Chapter 5, but with a different template file, `show_rpc.html`.

The displayed HTML, such as Figure 6.17, can be panned and zoomed via mouse drag and mouse wheel, respectively. The lower-left red dot resets the display. The [Rel. 0] button at the top switches between showing multiple RPCs by wall-clock- time and showing them all starting at (relative to) time zero.

6.16 Observations

On a sample server pair of machines, I measured 9.972 msec to send the Write() request from client to server, 1.118 msec to process it on the server, and 10 usec to send the response back to the client. My estimates of the first and last times were pretty reasonable, but my processing time estimate was 10x too low (see Table 6.1).

Table 6.1 Estimated RPC Elapsed Time vs. Measured

Action	Estimate	Actual
Send 1MB request	10 msec	9.972 msec
Process request	100 usec	1118 usec
Send 100 byte response	2 usec	10 usec

Based on the client-side log, Figure 6.17 shows 10 such RPCs lined up on time across about 120 msec, using the style of Figure 6.8. The hollow part of the lines (more precisely, the white overlaid lines) show approximately the message transmission time for the 1MB client => server messages. Among other things, you can see the 1–2 msec gap on the client from the end of one RPC to the beginning of the next one (oval around the first gap).

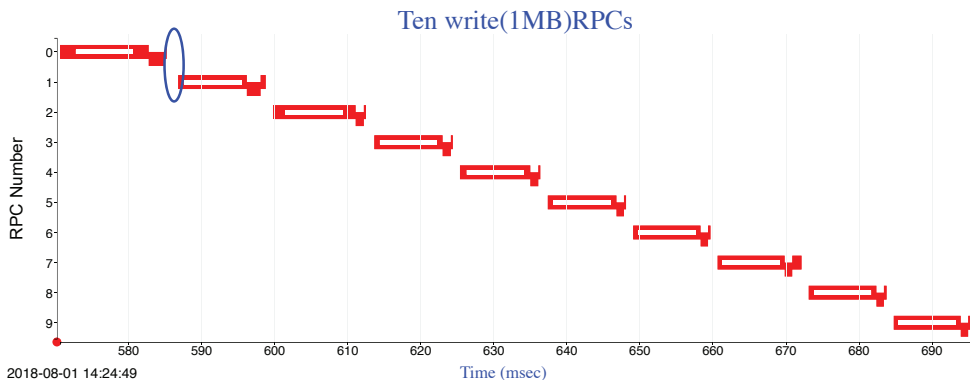


Figure 6.17 Ten RPCs, each sending 1MB of data. The oval highlights the 1–2 msec gap between successive RPCs.

Figure 6.18 shows the same 10 RPCs, but this time they are aligned to a common starting point, so the relative times of the different RPCs are easily compared. Now it is clearer that the first few RPCs take longer, and then the timing settles down a bit. It is also clearer that number seven has an extra delay in getting its response message back to the client.

The request time from client to server is about 9 msec, as we estimated, but the time processing on the server is about 1 msec, longer than the 100 usec per 1MB that we estimated. Why was our estimate off?

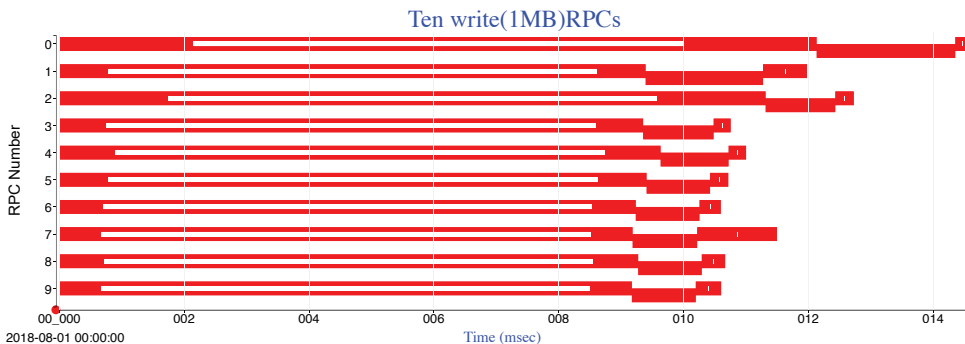


Figure 6.18 The same 10 RPCs, but this time all with the same start time of zero

Copying each 1MB message from network card to kernel buffer, then kernel buffer to user buffer, then user buffer to separate key and value strings, and then value string to map entry, each one MB in our sample server code is read or written to memory about eight times, not just once. This boosts our estimate to about 800 usec for all those copies, close to the observed 1 msec. To complete the picture, the response time from server to client is about 350 usec, a little longer than I would have expected but not terrible. Estimating ahead of time makes it easy to spot discrepancies.

6.17 Summary

In this chapter we introduced a sample remote procedure call database system and ran it on two networked computers, recording the times at which each RPC message is sent and received. The resulting notched-line diagrams show the measured time, the approximate transmission time for each message, the spacing between successive RPCs, and the relative timing of multiple similar RPCs. The underlying logging is sustainable at a rate of 10,000 RPCs per second per service.

In the next chapter, we will look in more detail at multiple overlapped RPCs on multiple clients and servers, at locking within transactions, and at time-aligning the clocks on multiple machines. We will expand the sample in-memory database to be an equally simple disk database. Later in the book we will examine the slow parts of these RPCs, including identifying the delays in message transmission and the delays in RPC processing.

- To understand the dynamics of multiple remote procedure calls, it is absolutely necessary to design in observation hooks, including at least RPC IDs, send and receive times, and byte lengths.
- These hooks must have low enough overhead to be useful under heavy live load.
- Creating RPC traffic requires at least one client and one server program running on different machines.
- A little stylized language lets us build a client to generate useful sequences of RPCs.

- Data structures for the RPC format across a network are part of the overall RPC design, and these must include observation metadata in the messages sent across a network, along with the actual operation and its data. This metadata information can fit in about 100 bytes.
- Logging to disk the timestamps and other metadata for every incoming and outgoing RPC message lets us observe a complete picture of where all the time went in an RPC call tree, and also a complete picture of all the *other* RPCs whose processing overlaps with and therefore possibly interferes with one of interest.
- Viewing the log data by wall-clock time lets us see delays within and between successive RPCs, and also lets us see overlapping of RPCs.
- Viewing the log data with all RPCs of a given kind starting at time zero lets us see differences among similar RPCs, and in particular lets us see what is different about slow ones vs. normal ones.
- Estimating what we expect to see makes it easy to spot discrepancies.

Exercises

Consider this work:

1. Send 10 ping messages of 100KB each.
2. Send 10 writes of 1MB of random data for keys kkkkk, kkkkl, kkkkm, ..., kkkkt.
3. Send 10 matching reads of 1MB from the same 10 keys.
4. Finally, send a quit command.

Draw yourself a little sketch of what you expect to see in the RPC timings.

Now run the `server4` program on one sample server and the `client4` program on another, sequentially sending commands for the previous sequences. Run the `dumplogfile4` program and `makeself` program against the first three client log files and display the actual results.

You likely will find that the two servers' wall-clock times differ by a few milliseconds, which may be enough to make the HTML display look odd, if the send time for a message is timestamped after the receipt time. We will look at time alignment in the next chapter. In the meantime, you might consider hand-editing the JSON files to adjust T2 and T3 to be between T1 and T4. This is optional, but doing so will give you some insight about what your Chapter 7 program will need to do.

- 6.1 How long, in milliseconds, did you estimate for the ping requests and their response message transmissions? How long do they actually take? Briefly comment on the difference.
- 6.2 How long, in milliseconds, did you estimate for the write requests and their response message transmissions? How long do they actually take? Briefly comment on the difference.
- 6.3 How long, in milliseconds, did you estimate for the read requests and their response message transmissions? How long do they actually take? Briefly comment on the difference.

This page intentionally left blank

Index

A

- Abusive offered load**, 162, 190
- ACK**, selective, 332, 333
- ACK**, TCP
 - n* Ethernet packets, 88
 - routing dynamics, 90
 - timeouts and, 358
 - transmission delays and, 326, 327, 331-333
- Acorn RISC machine**. See **ARM**
- AcquireSpin** loop, 343
- AcquireWait** routine, 342-344
- ACS-1** computer, 17-18
- ACS-360** computer, 18
- add** instruction
 - latency, 20-21
 - measurements, 23
- addevent** routine, 342
- addr2line** program, 264
- Address space layout randomization**. See **ASLR**
- Agarwal, Anant**, 195, 196, 405, 411
- Aggregate measures**, 141
 - conclusion, 155
 - histogram time scales, 147-149
 - measurement intervals, 143
 - patterns of values over time, 151-152
 - per-event, 150-151
 - timelines, 143-147
 - transaction examples, 154-155
 - uniform vs. bursty event rates, 142
 - update intervals, 152-154
- aio_select** service, 358
- Aligned** reference, 36
- Alignment**, 36, 111-117
- Allen, Frances E.**, xix, xxv
- Allen, K. Scott**, 338, 405
- Alto**, Xerox, 87
- AMD/Intel x86-64** patches, 236-237
- AMD Ryzen**, 387-388
- Amdahl, Gene M.**, 32, 345, 405
- Analog**, 64, 66, 359
- Analyzer programs** for observation tools, 170
- Aniszczyk, Chris**, 189, 405
- annot** group in software dynamics display, 258
- Annotations** in software dynamics display
 - timelines, 261-262
- Anomaly/ies**
 - echoing experiment, 334-335
 - at Google, 161
 - KUtrace for spotting, 225
 - reading from SSD, 80
 - resolution and spotting, 144
 - server4 program, 122
 - writing to SSD, 82
- Anonymous**, 146, 405
- Antagonist**
 - floating-point, 282-284
 - memory, 285
 - understanding, 286
- Apollo moon landing**, 358
- Arc** button in software dynamics display, 258, 260
- Arcs** in data display, 215
- arg0** field
 - JSON events, 256
 - trace entries, 231, 391-392
- ARM**, xxii, 173, 224, 233, 384
- Arms** in hard disks, 62-63
- Arrival times** on master dashboards, 159-160
- ASLR (address space layout randomization)**, 264
- Aspect** box in software dynamics display, 266
- Associativity** measurements for cache levels, 45-46
- Asynchronous**
 - effect of blktrace on, 185
 - defined, 397
 - read, 69, 73, 74
 - RPCs, 6, 86
 - with SSD, 80, 82
 - write, 72, 77, 78
- Atomic instructions**
 - locks, 338, 344
 - queues, 365, 376
- Atomic pointer increments**, 241-242
- ATIME** directory field, 67, 397
- Atlas**, Manchester, 17, 412
- Averages**
 - aggregate measures, 145-147
 - CPU-only performance, 168
 - latency, 4
 - top command, 170
- Axes**
 - data display, 215
 - software dynamics display, 265
- axisLabelX** field for JSON metadata, 254
- axisLabelY** field for JSON metadata, 255

B

- Babcock, Charles, 136, 405
 - Back end, execution, 6, 7, 19
 - Back-of-the-envelope calculation, 71, 99, 158, 215
 - Bakita, Joshua J., xxv
 - Balan, Subra, 409
 - Bannon, Pete, xxv
 - Barroso, Luiz André, xix, 411
 - Batch software, 4, 5, 171, 327
 - Battle, T. R., 407
 - Beadon, Matt, 407
 - Beaver, Donald, 411
 - Benchmarks
 - floating-point antagonists, 283–284
 - profiles, 173–176
 - slow performance problem, 279–282
 - SPEC, 58
 - Bignose, 406
 - Binary log file records, 137–138
 - Binary order of magnitude, 9
 - Biro, Ross, 230
 - Blaauw, Gerrit A., 338, 405
 - Blake, Geoffrey, 8, 405
 - Blind to, 168, 193, 210, 370, 385
 - blktrace** command, 184–187
 - Boatload, 5
 - Boggs, David R., 87, 414
 - Bond, Elaine, xxv
 - Borg, Anita, xxv, 211, 406
 - Bostock, Mike, 249, 406
 - Bowen, M. Lucey, xxv
 - Branch instruction latency, 20
 - Brockmeyer, R. L., 407
 - Brooks, Frederick P., Jr., xix, xxv, 338, 405, 406
 - Broughton, Jason, 327, 406
 - Brown, Michael, xxv
 - Buffers
 - disks, 66–68
 - early CPUs, 17
 - KUtrace tool, 223
 - Linux kernel patches, 227
 - log files, 138
 - organization, 36–37
 - trace, 205
 - trace buffer data structures, 228
 - BuildDashString** routine, 354
 - Burrows, Mike, 411
 - Bursty events, 142, 146, 370, 397
 - Busy machines, 222, 223, 315, 317
 - Byte addressing, 32, 36
-
- C**
 - C6 sleep state, 262, 293, 294
 - c_str()** routine, 184
 - Caches
 - AMD Ryzen 3 chip, 387–388
 - associativity, 37, 38, 45, 71
 - cache-aware computation, 58
 - cache hits, 35, 38, 44, 45, 51
 - cache size, 37, 40, 44, 46, 47
 - CPU and memory interaction, 49–51
 - datacenters, 6
 - direct-mapped, 36, 37
 - disks, 67
 - hashed lookup, 54, 55, 56, 57
 - hierarchy, 31, 34, 35, 37, 43, 400
 - introduction, 18
 - KUtrace tool design, 221
 - level associativity measurements, 45–46
 - level size measurements, 43–45
 - lines, 35–39, 41–51, 211, 221
 - memory, 18, 33, 34–35, 400
 - miss, 6, 35, 38, 39, 50–55, 71, 201
 - organization, 34–36
 - physically addressed, 37, 47
 - prefetching, 40–42, 47, 57, 61, 83
 - replacement, 44, 400, 401, 402
 - set-associative, 36, 45, 49, 402
 - sets, 50, 51, 54, 57
 - tags, 35, 36, 37, 403
 - total size, 38, 43, 44, 45
 - underutilization, 44, 46
 - virtually addressed, 37
 - Calculation intervals on master dashboards, 160
 - Canary requests in network traffic, 95
 - CB** button for software dynamics display, 259
 - CBS, 358, 406
 - CDC 6600 computer, 17
 - CFS (Completely Fair Scheduler), 289–292
 - Chakravarthy, Tejasvi, 407
 - Checksums
 - packets, 87–88
 - RPCs, 98, 100, 104
 - Chesson, Greg, 196, 406
 - ChipsEtc, 33, 406
 - chksum** method
 - RPC clients, 104
 - RPC logs, 100
 - Christie's, 358, 406
 - Clark, Douglas, 407
 - Client-mode RPCs, 92
 - Client-server systems, 7, 100, 105–106, 397
 - Clients and **client4** program
 - disk and network database interaction, 117–125
 - latency in master dashboards, 161
 - network waiting issues, 320–321, 327–329
 - RPCs, 100, 102–106
 - time-of-day clock, 113
 - Clock cycle, 15, 17, 25, 32
 - Clock drift, 112–114
 - Clock interrupts, 168
 - Clock skew, 111
 - Cocke, John, xix, xxv
 - Colossus file system, 408

- Collector programs for observation tools, 170
- Color in data display, 215
- Column access array, 50
- Column access, DRAM, 42
- Comcast, 189, 406
- Comment field for JSON metadata, 254
- Commitments for service response times, 190–191
- Common Vulnerabilities and Exposures. *See* CVE
- Compiler optimization, 24, 28, 58, 284
- Completely Fair Scheduler (CFS). *See* CFS
- Compression in SSDs, 65
- Constant-rate timestamp counters, 21
- Constraints
 - offered loads, 11
 - trace tools, 202
- Contented locks
 - critical section locks, 189, 196
 - waiting for. *See* Locks, waiting for
- Contented-wait-time histograms, 339
- Context switches in disk operations, 309
- Conti, Carl J., 36, 406
- Control interfaces for loadable modules, 247
- Control programs for KUtace tool, 223–224
- Controls for software dynamics display, 258–259, 265–266
- Coordinated universal time (UTC) for timestamps, 136
- `copy_from_user()` routine, 244
- Copy-on-write (CoW). *See* CoW
- Copying
 - disk data, 78–80
 - log files, 136
 - SSD data, 82
- Core memory
 - CPU access, 15–16
 - cycle times, 33
 - ferrite, 32
- Counters
 - logs, 133
 - observation tools design, 212
 - oprofile, 173
 - overview, 167–169
 - perf, 171–173
 - time, 171
 - top, 170
- Counting instruction cycles, 20–22
- CoW (copy-on-write), 293
- CPU affinity, 294, 295, 296, 297
- CPU and memory interaction, 49
 - cache-aware computation, 58
 - cache interaction, 49–51
 - exercises, 59
 - matrix multiplication algorithm, 51
 - matrix multiplication estimates, 51–52
 - matrix multiplication initialization,
 - cross-checking, and observing, 52–53
 - matrix multiplication subblock method, 57–58
 - matrix multiplication transpose method, 55–57
 - summary, 58–59
- CPU architecture suggestion, 37, 201, 232
- cpu field for JSON events, 255
- CPU group in software dynamics display, 259–260
- cpuModelName field for JSON metadata, 255
- CPUs
 - aggregate measures, 151
 - frequency, 20, 259, 270, 369, 370, 384
 - IPC trace entries, 232
 - KUtace tool design, 222
 - lock waiting issues, 342
 - measuring. *See* CPUs, measuring
 - memory hierarchy, 34–36
 - memory interaction. *See* CPU and memory interaction
 - oprofile profiler, 173–176
 - per-function Gmail trace example, 203–207
 - per-server dashboard information, 164
 - profiling tools, 168
 - queue issues, frequency, 369–370
 - queue issues, KUtace, 371–374
 - queue issues, RPC logs, 370–371
 - software dynamics display, 259–260
 - software dynamics display timelines, 260–261
 - tool overhead, 211–212
 - top command, 170
 - trace buffer data structures, 228
 - trace tools, 195
 - waiting for. *See* CPUs, waiting for
- CPUs, measuring
 - add instruction latency, 20–21
 - CPU history, 15–18
 - dead variable fail, 24–25
 - dependent variables, 26
 - execution latency, 26–27
 - exercises, 28–29
 - `gcc -fno-tree-reassoc` command-line flag, 27
 - loop fails, 21–24
 - loops with volatiles, 25–26
 - modern processor chips, 19
 - straight-line code fail, 21
 - summary, 28
- CPUs, waiting for
 - analysis, 293–295
 - exercises, 297
 - exploring and reasoning, 290–292
 - mystery, 289–290, 292
 - program, 289
 - scheduler issues, 289–292
 - software locks, 295–297
 - startup issues, 292
 - summary, 297
- Cray-1 computer, 20

Critical sections

- description, 12
- lockttrace command, 189
- queue issues, 375–377
- spinlocks, 101–102
- trace tools, 196
- waiting for. *See* Locks, waiting for

Cross-checking

- matrix multiplication, 52–53, 55
- queue issues, 379–380

Crystal oscillators for clock, 111–112**ctime() routine, 137****Culler, David E., 408****Curved arcs in data display, 215****CVE 2017, 225, 406****Cycles**

- counting, 20–22
- CPUs, 15–17
- SSD erase, 65

Cylinder, disk, 62–64, 69, 70, 397**D****D-cache, 18, 387, 388****D3 JavaScript library, 249, 254****Damato, Joe, 327, 406****Dapper tool**

- traffic observations, 94
- working with, 189–191

dashboard_thread processes, 341–344**Dashboards, 157**

- BirthdayPic sample service, 157–158
- data display, 214
- exercises, 165–166
- lock issues, 341–344, 353–354
- master, 159–163
- per-instance, 163–164
- per-server, 164
- sample, 159
- sanity checks, 164–165
- summary, 165

Data alignment in memory, 36**Data display**

- observation tools design, 214–215
- software dynamics display. *See* Display of software dynamics

Data-driven documents. *See* D3**Databases**

- datacenter software layer, 6–7
- disk and network database interaction. *See* Disk and network database interaction
- sample server program, 101

Datacenters

- hardware, 5–6
- software, 6–7
- terms and concepts, 3–5

datalen field in RPC headers, 98**Dates in data display, 215****De Bruijn, Willem, 408****Dead code, 24, 28****Dead variable fail measurements, 24–25****Dean, Jeffery A., 8, 10, 95, 406****DEC Alpha 21064, 18, 20****Decimal orders of magnitude, 9****delete method**

- RPC clients, 104
- RPC logs, 100

Delta time in disk and network database interaction time, 113–115**Deming, W. Edwards, 1****Dennard, Robert H., 33, 406, 412****Dependent instructions**

- defined, 398
- loads and memory, 41
- variables, 26

Dequeue routine, 365–366, 374–375**Design**

- KUtrace tool, 221–223
- observation tools. *See* Observation tools design
- traces, 194–197

Destructive readout, 33, 142, 398**Detour, intentional. *See* Enticement****Digital Equipment Corporation. *See* DEC****DIMMs (dual-inline memory modules), 32, 42****Direct-mapped caches, 36****Directory, file system, 66, 67, 121, 171****Disaster, performance, 4, 67, 68, 201, 202****Disk and network database interaction, 111**

- exercises, 129
- experiment 1, 118–121
- experiment 2, 121–125
- experiment 3, 125–127
- logging, 127
- multiple clients, 117
- on-disk databases, 121
- spinlocks, 118
- summary, 128–129
- time alignment, 111–117
- transaction latency variation, 128

Disks

- disk blocks, 63, 71, 255, 402
- disk syncs, 308–310, 314–316
- network interaction. *See* Disk and network database interaction
- per-server dashboard information, 164
- tool overhead, 211–212
- trace tools, 195
- tracing with blktrace command, 184–187
- waiting for. *See* Disks, waiting for

Disks/SSD, measuring, 61

- back-of-the-envelope calculations, 71–72
- disk reads, 68–71, 73–77
- disk writes, 77–80
- exercises, 84

- hard disk overview, 62–64
- multiple transfers, 82–83
- software disk access and on-disk buffering, 66–68
- SSD overview, 64–66
- SSD reads, 80–81
- SSD writes, 82
- summary, 83
- Disks, waiting for, 307**
 - analysis, 317
 - exercises, 317
 - exploring and reasoning, 308–310
 - mystery, 307–308
 - program, 307
 - read and write times, 307–309
 - reading 40MB on disks, 310–311
 - reading 40MB on SSD, 315
 - reading random 4KB blocks, 313–314
 - reading sequential 4K blocks, 311–313
 - summary, 317
 - two programs accessing two files at once, 316–317
 - writing and sync of 40MB on SSD, 314–315
- Display, observation tools design for, 214–215**
- Display of software dynamics, 257**
 - controls, 258–259
 - IPC legends, 265
 - overview, 257–258
 - Save/Restore, 265
 - secondary controls, 265–266
 - summary, 266
 - timelines, 260–265
 - X-axis, 265
 - Y-axis, 259–260
- Distortion, 21, 25, 42, 133, 148, 194**
- Diwan, Amer, xxv, 353**
- Dixit, Harish Dattatraya, 211, 407**
- Dixit, Kaivalya M., 58, 406**
- Doench, Greg, xxv**
- DoFakeWork routine, 342**
- DRAM (dynamic random-access memory)**
 - non-random, 42–43
 - overview, 33
 - sample server, 32
- dt field in trace entries, 231, 391**
- Dual-inline memory module. See DIMM**
- dumplogfile4 program, 100, 106**
- duration field for JSON events, 255**
- Dwarkadas, Sandhya, xxv, 222**
- Dynamic random-access memory. See DRAM**
- Dynamics of programs, 4**

E

- E field in trace entries, 231, 391**
- ECC (error correction code) bits**
 - memory, 33
 - SSDs, 66
- Einstein, Albert, 131**

- Elapsed time in data display, 215**
- Embedded servo, disk, 63, 64, 75, 398**
- Emer, Joel S., xxv, 196, 407**
- Empirical transactions per second targets, 4**
- Encapsulated packets, 91**
- Encrypted packets, 91**
- Enqueue routine, 366, 374–376**
- Enticement, seductive. See Mislead**
- Erase cycles, SSD, 65, 68, 398**
- Error correction code. See ECC**
- Errors per interval data on per-instance dashboards, 163**
- Estimates, 9–10, 51**
- Ethernet, 87–88**
- Ethertype, 88, 188, 321**
- ethtool, 327**
- Event numbers**
 - KUtrace tool, 219
 - Linux kernel patches, 233
 - trace entries, 391–395
- Events**
 - JSON format, 255–256
 - uniform rates vs. bursty, 142
- eventspan program, 251–252, 262–263**
- Execute cycle, 16, 20**
- Execution latency in CPU measurements, 26–27**
- Execution units, 16, 19**
- Extended logging, 135**
- Extent, file, 66–67, 69–71, 77, 398**
- Extinction event, 143–145, 411, 413**
- Extra work, 11, 277, 358**
- Eyles, Don, 358, 407**

F

- Fade button for software dynamics display, 266**
- Fall, Kevin R., 321, 407**
- FancyLock class, 339**
- Farrow, Rik, xxv**
- fdiv_wait_usec routine, 367**
- Ferrite core memory, 32**
- Fetch/execute CPU process, 16**
- Feynman, Richard P., 267**
- Files**
 - log formats, 137–138
 - metadata, 67
- Fills in cache lines, 36**
- Filters in packet tracing, 235**
- FindCacheSizes() routine, 44**
- Fixed-length trace entries, 391–392**
- flags field for JSON metadata, 255**
- Flash drive. See SSD**
- Flash memory, 64–66**
- Flight-recorder mode**
 - KUtrace tool, 221
 - observation tools, 170
 - postprocessing, 250
 - trace calls, 243
 - trace tools, 197

Floating gates in SSDs, 64, 66
 Floating-point antagonists in slow performance problem, 282–284
flt_hog program, 282–284
 Fogh, Anders, 409
 Formats for log files, 137–138
 Four-way set-associative caches, 36
fprintf routine, 354
 Fram oil filter, 124, 407
free routine, 183–184
freq button for software dynamics display, 259–260
 Frequency decisions in observation tools design, 210–211
 Friedenber, S. E., 196, 197, 407
 Friend of the Electron, v
 Friend, J. R., 407
 Front end, instruction fetch, 19
ftrace command, 180–183
 Full intervals for aggregate measures, 154
 Fully associative caches, 36
 Fundamental resource, 12, 15, 31, 169, 195, 262
futex routine
 CPU waiting issues, 295–296
 spinlocks, 376–377
futex_wait routine
 lock issues, 343–344
 queue issues, 373, 375
futex_wake routine, 373, 375
fwrite() routine, 137

G

Gallatin, Drew, xxv
 Galloping Gertie, 413
 Garden path, 379; *See also* Detour
gcc -fno-tree-reassoc command-line flag, 27
 GE 645, 17
 Generators for datacenters, 5
 Geng, Yilong, 113, 407
 Genkin, Daniel, 409
get_claim routine, 241–244
get_slow_claim() routine, 244
gettimeofday() routine
 disk reads, 68
 disk writes, 71
 locks, 343
 log files, 137
 network waiting issues, 324–325
 postprocessing, 250
 time-of-day clock, 112
 timestamps, 135–136
 traceblocks, 230
GetUsec() routine
 disk reads, 68
 lock waiting, 343
 GFS (Google File System), 125, 407
 Ghemawat, Sanjay, 125, 407
 Gifford, David, 338, 407

Girl with a curl, 31, 409
 Glitch, hardware, 210
 Glossary of terms, 397–404
 Gmail per-function trace example, 203–207
go command in **kutrace_control** program, 246
 Goals for KUTrace tool, 220–221
 Gobioff, Howard, 407
 Goldberg, Patricia A., xxv
 Golden age of computing, 15, 32, 399
 Goldilocks, 220
 Goldsmith, Belinda, 146, 407
 Good learning, 361, 379, 386
 Good trouble, 361, 409
 Google, 138, 189, 196, 407
 Google File System. *See* GFS
 Govindaraju, Rama, 408
 GPU (Graphics Processing Units), 220, 385
 Gregg, Brendan D., 11, 167, 407
 Gross interference data in per-server dashboards, 164
 Gruss, Daniel, 409
 Guaranteed-delivery protocols, 88, 90
 Gustafson, John L., 345, 408

H

Haas, Werner, 409
 Half-Optimal Principle
 CPU waiting problem, 293, 295
 SSD reads, 315
 Half-Useful Principle
 context switching, 309
 disk reads, 314
 disk writes, 71–73
 latency, 69–70
 SSD reads, 315
 Hamburg, Mike, 409
 Hard disks. *See* Disks
 Hardware
 datacenters, 5–6
 performance counters, 168
 shared, 12
 Hardware interference with KUTrace tool, 222
 Hashes for memory, 36
hdparam -W0 command, 67
 headerlen field in RPC headers, 97
 Headers
 Ethernet, 88
 JSON, 252
 packets, 90–91
 RPC messages, 97–99, 263
 Heads
 hard disks, 62–64
 queue structure, 364
 Head-switching, disk, 64, 75
 Heat maps, 151–152
 Hedberg, Ray, xxv
 Heisenberg, Werner, 194
 Held locks, 337–338

- Hennessy, John L., 15, 408
 Herbert, Tom, 335, 408
 High-density data display, 215
 Hildebrand, Dean, 125, 408
 Hindered transactions, 11, 13, 399
 Histograms
 aggregate measures, 145–146
 case study, 212–214
 latency, 8–9
 lock waiting, 339
 long-tail distributions, 145–146
 per-instance dashboards, 163–164
 small, 102, 189
 spinlocks, 102
 time scales, 147–149
 Hits
 cache, 35
 description, 6
 Hochschild, Peter H., 211, 408
 Hoff, Todd, 8, 408
 Hog, floating-point, 280, 282–286
 Hog, memory, 280,
 Hog, paging, 299–304, 306
 Horn, Jann, 409
 Horowitz, Mark, 405
htop command, 170
 Hubs, 88, 89
 Hunt, V. Bruce, xxv, 410
 HyperThread, 18
 Hypervisors, 385
-
- I**
- I-cache, 18, 172, 387, 388
 IBM 1959, 199, 408
 IBM 1967, 338, 408
 IBM 1970, 33, 408
 IBM 1983, 376, 408
 IBM 2021, 189, 408
 IBM 33FD, 186
 IBM 360/67, 17
 IBM 360/85, 18, 33, 34, 35
 IBM 360/91, 411
 IBM 370/145, 408
 IBM 704, 16
 IBM 709, 16
 IBM 7010, 197–199
 IBM 7030, 16, 412
 IBM 7094, 16
 IBM Power4 chip, 18
 IBM RISC System/6000, 18
 IBM Personal Computer AT, 65
 IBM Power4, 18, 412
 IBM Stretch, 16–17
 IBM System/360 computer test-and-set
 instruction, 338
 IBM System 370/145 computer, 33
 Idle, low-power, 262, 313, 384, 394
 IEEE 2021, 88, 408
 Information Sciences Institute, 333, 408
 Initialization
 matrix multiplication, 52–53
 tracing, 241
 Inline execution delays, 359
 Input/output memory management unit.
 See IOMMU
 Insert1 routine, 241–243
 InsertN routine, 243–244
 Instruction fetch, 15, 16, 18, 21, 25, 32, 198, 283
 Instruction sets, 16–17
 Instructions executed per cycle (IPC)
 KUtrace tool design, 222
 legends in software dynamics display, 265
 trace buffer data structures, 228
 trace entries, 232
 Intel, 33, 232, 236, 293, 327
 Intel 1103, 33
 Intel 2021, 172, 408
 Intel i3, 32, 38, 48, 281, 289, 295, 388
 Intel Pentium 4 processor, 18
 Intel Pentium P5 processor, 20
 Intel Xeon, 18
 Interactive data displays, 215
 Interface data structures, 239–240
 Interference, 143, 150, 195, 217, 222, 280,
 368–369
 Intermittent performance problem
 exploring and reasoning, 273–276
 mystery, 272–273
 mystery21 program, 271–277
 overview, 271
 summary, 277–278
 understanding, 277
 Internet protocol. See IP
 Interrupts
 coalescing, 325, 327
 delays, 329
 disk operations, 308–313
 network waiting, 327
 profiling tools, 168
 request. See IRQ
 Intervals
 aggregate measures, 143, 152–154
 master dashboards, 160–162
 IOMMU, 82, 399
 IP addresses in TCP/IP, 90
 IP field in RPC headers, 98
 IPC. See Instructions executed per cycle (IPC)
 IPC button for software dynamics display, 259
 ipc field for JSON events, 256
 IPv4, 88, 90, 235, 321, 399
 IPv6, 88, 235, 321, 399
 IRQ, 224, 234, 236, 327, 392
 Issue cycle, 343, 399
 Issue slots, 19

J

Jaspan, Saul, 411
 Jeffay, Kevin, xxv
 Jiffie, 325
 JSON format
 events, 255–256
 metadata, 254–256

K

Karp, Joel, 33, 412
 Kaufmann, Richard, xxv
 Keatts, Bill, 58, 408
 Keller, James B., xxv
 Kernel buffers for KUtrace tool, 223
 Kernel interface data structures, 239–240
 Kernel-mode code for RPC messages, 320
 Kernel-mode samples for profiling tools, 168
 Kernel patches and module
 KUtrace tool, 224, 394
 Linux. *See* Linux kernel patches for KUtrace
 Kernel-user trace. *See* KUtrace tool
 kernelVersion field for JSON metadata, 255
 Kernighan, Brian, xxv
 Kessler, Richard E., 406
 Knuth, Donald E., xix, xxv, 19, 196, 409
 Koher, Paul C., 225, 409
 Kozyrakis, Christos, xxv
 Krohnke, Duane W., 17, 409
 kswapd0 daemon, 300–301
 ktime_get_real routine, 324–325
 Kuck Associates, 58
 kutrace_control program, 246–247
 kutrace_global_ops global, 234, 240
 kutrace_lib library, 246–247, 393
 kutrace_mod_exit routine, 240
 kutrace_mod_init routine, 240
 kutrace_mod module, 264
 KUtrace tool, 219
 control program, 223–224
 CPU waiting problem, 290–292
 design, 221–223
 goals, 220–221
 implementation, 223
 JSON format, 254–256
 kernel patches and module, 224
 Linux kernel patches. *See* Linux kernel patches for KUtrace
 Linux loadable module. *See* Linux loadable modules for KUtrace
 overview, 219
 postprocessing. *See* Postprocessing for KUtrace
 queue waiting, 371–374
 security, 225
 software dynamics display. *See* Display of software dynamics
 trace entries, 391–395

user-mode runtime control. *See* User-mode runtime control for KUtrace

kutrace_tracing global, 234, 240
 kworker0:0 process, 301

L

L fields in RPC headers, 98
 L1 cache, 35, 37, 43–46, 49–51, 294–295
 L2 cache, 18, 34, 44, 201, 294–295
 L3 cache, 34, 44, 49, 51, 54–55, 295
 Labels in data display, 215
 Lamport, Leslie B., 242, 409
 Larus, James R., xxv
 Latency
 add instruction, 20–21
 aggregate measures, 147
 description, 3
 disk. *See* Disks/SSD, measuring
 disk and network database interaction, 128
 execution, 26–27
 Half-Useful Principle, 69–70
 instruction, 28, 38, 399
 light load, 161
 load-to-use, 41, 44
 long tail, 7–9, 11, 147–149, 267, 400
 master dashboard times, 160–161
 memory, 31
 queries, 4
 queues, 368–369
 response, 161
 RPC, 92
 tail, 4, 9, 222
 tail-flick, 403
 transaction, 189–191
 Layers
 datacenter subsystems, 6–7
 memory access, 31–32
 LBAs (Logical Block Addresses), 63, 65, 66, 399
 Leap seconds with timestamps, 136
 Least recently used. *See* LRU
 Lee, Ki Suh, 113, 409
 Legend button for software dynamics display, 266
 Legends in software dynamics display, 265
 Length decisions in observation tools design, 210–211
 Leung, Shun-Tak, 407
 Level-1 instruction caches, 34–35
 Levels, cache
 associativity measurements, 45–46
 size measurements, 43–45
 Lewis, John R., 361, 409
 Line size for caches, 38–39
 Linear feedback shift registers, 137
 Linear scales for aggregate measures, 147–148
 LinearTiming() routine, 41
 Lines in cache memory, 35
 Linux, 196, 289, 409

- Linux kernel patches for KUtrace, 227**
 - AMD/Intel x86-64 patches, 236–237
 - code, 234
 - event numbers, 233
 - exercises, 237
 - IPC trace entries, 232
 - nested trace entries, 233
 - packet tracing, 234–235
 - raw traceblock format, 229–230
 - summary, 237
 - syscall trace entries, 230–231
 - timestamps, 233
 - trace buffer data structures, 228
 - Linux loadable modules for KUtrace, 239**
 - kernel interface data structures, 239–240
 - loading/unloading, 240
 - summary, 244
 - trace calls, implementing, 241
 - trace calls, initializing and controlling, 241
 - trace calls, Insert1 routine, 241–243
 - trace calls, InsertN routine, 243–244
 - traceblock switching, 244
 - Lipp, Moritz, 225, 409**
 - Liptay, John S., 18, 33, 409**
 - list command, 171–172**
 - Liu, Shiyu, 407**
 - Load balancing for queues, 377–378**
 - Loadable modules**
 - control interfaces, 247
 - Linux. *See* Linux loadable modules for KUtrace
 - Lock button for software dynamics display, 259**
 - Lock capture, 338, 347–348**
 - Lock contention, 349–354**
 - Lock saturation**
 - defined, 338
 - long lock hold times experiment, 345
 - Locks**
 - CPU waiting issues, 295–297
 - disk and network database interaction, 118
 - software dynamics display timelines, 263
 - spinlocks, 101–102
 - tools, 169
 - Locks, waiting for, 337**
 - experiment 1, exploring and reasoning, 346–348
 - experiment 1, lock capture, 347–348
 - experiment 1, lock saturation, 345
 - experiment 1, lock starvation, 348
 - experiment 1, mysteries, 345–346
 - experiment 1, simple locking, 344
 - experiment 2 lock capture fixes, 348–349
 - experiment 3 lock contention fixes, 349–351
 - experiment 4 lock contention fixes, 351–353
 - experiment 5 lock contention fixes, 353–354
 - overview, 337–341
 - program, 341–344
 - summary, 355
 - locktrace command, 189**
 - Logarithmic scales for aggregate measures, 147–148**
 - Logging and logs**
 - basic, 134–135
 - disk and network database interaction, 127
 - extended, 135
 - file formats, 137–138
 - managing, 138–139
 - networks, 99–100
 - observation tools, 133, 212
 - overview, 133–134
 - postprocessing, 106
 - queue waiting, 370–371
 - RPC IDs, 136–137
 - RPCs, 100
 - summary, 139
 - timestamps, 135–136
 - tracing tools, 168
 - Logical Block Addresses. *See* LBAs**
 - Logical Block Numbers, 66–67**
 - Long lock hold times experiment**
 - exploring and reasoning, 346–348
 - lock capture, 347–348
 - lock saturation, 345
 - lock starvation, 348
 - mysteries, 345–346
 - simple locking, 344
 - Long-tail latency, 11, 267, 400**
 - histogram time scales, 147–149
 - overview, 7–9
 - Long-tail transactions with KUtrace tool, 222**
 - Longbottom, Roy, 279, 409**
 - Longfellow, Henry Wadsworth, 31, 409**
 - Loop fail measurements, 21–24**
 - Loop unrolling, 25, 28, 58**
 - Low overhead**
 - contended lock tools, 189
 - counting tools, 168
 - KUtrace tool, 217, 223, 227
 - network trace, 188
 - observation tools, 133
 - profile tools, 168, 176
 - protocol buffers, 138
 - RPC logs, 189
 - timestamp logging, 127
 - of top command, 170
 - tracing tools, 180, 193, 194, 195
 - Low-power, idle, 262, 384**
 - LRU, 44, 400**
 - ltrace command, 179–180**
-
- M**
- MAC (Media Access Control) addresses**
 - Ethernet packets, 88
 - switches, 89
 - Make-runnable changes in KUtrace tool design, 222**
 - MakeAction routine, 342, 346**

- MakeLongList() routine**, 42
- makeself program**, 254
- malloc routine**, 183–184, 299
- Manchester Atlas machine**, 17
- Mangard, Stefan**, 409
- Mapping, address**, 17, 37, 43, 47, 48, 52
- Mark button for software dynamics display**, 258, 260
- Markers in KUtrace tool**, 219
- Mason, Chris**, 407
- Massachusetts Institute of Technology**. *See* MIT
- Master dashboards**, 159–163
- matrix.cc program**, 53
- Matrix multiplication**
 - algorithm, 51
 - estimates, 51–52
 - initialization, cross-checking, and observing, 52–53
 - subblock method, 57–58
 - transpose method, 55–57
- Matrix storage order**, 49–51
- Matrix300 benchmark**, 58, 280
- Maurer, James**, xxv
- Maximum values in aggregate measures**, 147
- mbit_sec field for JSON metadata**, 255
- McKenney, Paul E.**, 354, 409
- Media Access Control addresses**. *See* MAC
- Measurement intervals**
 - aggregate measures, 143, 152–154
 - master dashboards, 161–162
- Measuring**
 - CPUs. *See* CPUs, measuring
 - disks/SSD. *See* Disks/SSD, measuring memory
 - memory. *See* Memory, measuring
 - networks. *See* Networks, measuring
- Medians**
 - aggregate measures, 145–147
 - latencies, 8
 - master dashboards, 162
- memcpy() routine**, 244
- Memory**
 - banks, 16, 19
 - CPU access, 16
 - CPU interaction. *See* CPU and memory interaction
 - hierarchy, 32, 34, 35, 37, 38, 46
 - Linux kernel patches, 227
 - measuring. *See* Memory, measuring
 - mttrace command, 183–184
 - protection, 17, 177
 - SSDs, 64–66
 - tool overhead, 211–212
 - top command, 170
 - trace tools, 195
 - waiting for. *See* Memory, waiting for word, 33, 34
- Memory, measuring**
 - cache level associativity, 45–46
 - cache level size, 43–45
 - cache line size, 38–39
 - cache organization, 34–36
 - cache underutilization, 46
 - data alignment, 36
 - dependent loads, 41
 - exercises, 47–48
 - history, 32–34
 - N+1 prefetching problem, 40–41
 - non-random DRAM, 42–43
 - process, 37–38
 - summary, 46–47
 - timing, 31–32
 - translation buffer time, 46
 - translation lookaside buffer organization, 36–37
- Memory, waiting for**, 299
 - analysis, 304–305
 - exercises, 306
 - exploring and reasoning, 300–304
 - mystery, 300, 304
 - program, 299
 - summary, 306
- Memory antagonists in slow performance problem**, 285
- Memory hog**, 280, 285
- Memory used data in per-server dashboards**, 164
- Metadata**
 - files, 67
 - JSON format, 254–255
- Metcalfe, Robert M.**, 87, 409, 414
- Method field in RPC headers**, 99
- Microprocessor superscalar design**, 18
- Minimum values in aggregate measures**, 147
- Mislead or deceive**. *See* Primrose path
- Misses**, 6, 35
- MIT (Massachusetts Institute of Technology)**, 17
- MLC (multi-level cell) drives**, 66
- Modern processor chip speedup techniques**, 19
- Modules for KUtrace tool**, 224
- Mogul, Jeffrey C.**, 408
- Monolithic memory**, 33
- mttrace command**, 183–184
- Mudge, J. Craig**, xxv
- Multi-bit error correcting codes in SSDs**, 66
- Multi-level cell (MLC) drives**. *See* MLC
- Multics operating system**, 17
- Multiple clients in disk and network database interaction**, 117–125
- Multiple disk transfers**, 82–83
- Multiple instruction issue**, 19
- Multiple threads**
 - code locking out, 214
 - global counter disaster and, 202

- in datacenter environment, 5, 12
- lock saturation and, 338, 344, 400
- RPC overlap and, 118
- RPCs handled by, 101
- spawned by PID 1234, 264
- Multiple time groupings in master dashboards, 161–162**
- Multiplication of matrices. See Matrix multiplication**
- multiply instruction, 16**
- Murray, Hal, xxv, 334**
- Mutex class, 339**
- Muthiah, Bharath, 407**
- mwait instruction, 293–294**
- mystery0 program, 22–26**
- mystery2 program, 38, 40–42, 44**
- mystery3 program, 73**
- mystery21 program, 271–277**
- mystery23 program, 289–291**
- mystery25 program, 308–312, 315–316**
- mystery27 program, 341–354**

N

- N+1 prefetching problem, 40–41**
- Naik, Ashish, 407**
- NaiveTiming() routine, 40–41**
- name field for JSON events, 256**
- Names for KUltrace tool intervals, 221**
- nanosleep routine, 327–328**
- National Institute of Standards and Technology. See NIST**
- Naur, Peter, 410**
- Nested trace entries, 233**
- Network interface controller. See NIC**
- Network time protocol. See NTP**
- Networks**
 - bandwidth data on per-server dashboards, 164
 - disk interaction. *See* Disk and network database interaction
 - measurements. *See* Networks, measuring
 - packets. *See* Packets
 - tools, 187–189
 - trace tools, 195
 - waiting for. *See* Networks, waiting for
- Networks, measuring, 85–87**
 - client program, 102–105
 - Ethernet, 87–88
 - exercises, 109
 - hubs, switches, and routers, 89
 - logging design, 99–100
 - observations, 107–108
 - packets, 90–91
 - RPCs, client/server systems, 100, 105–106
 - RPCs, log postprocessing, 106
 - RPCs, message definition, 96–99
 - RPCs, overview, 91–93
 - server program, 101
 - slop, 93–94
 - spinlocks, 101–102

- summary, 108–109
- TCP/IP, 89–90
- traffic observations, 94–96
- Networks, waiting for, 319**
 - analyses, 333
 - anomaly, 334–335
 - experiment 1, exploring and reasoning, 323–327
 - experiment 1, mystery, 322–323
 - experiment 1, overview, 321
 - experiment 1, time between RPCs, 327–329
 - experiment 2, 329
 - experiment 3, 329–330
 - experiment 4, 330–333
 - overview, 319–320
 - programs, 320–321
 - summary, 336
- Neural network processing chips, 385**
- Nguyen, Thomas, 413**
- NIC (Network Interface Controller), 324**
- NIST (National Institute of Standards and Technology), 136, 410**
- Non-blocking RPC, 96–97**
- Non-execution spans, 262–263**
- Non-overlapping intervals in aggregate measures, 152–153**
- Non-random DRAM, 42–43**
- Nop events in trace entries, 391**
- Normal distributions for aggregate measures, 145–147**
- Not pipelined, 367, 369**
- Notches in data display, 215**
- Nothing missing design, 190, 193, 205, 207, 216, 220, 225, 385**
- NTP (Network Time Protocol), 112**
- Nyland, Lars, 289**

O

- O_DIRECT disk-read pattern, 310**
- O_DIRECT parameter, 67**
- O_NOATIME parameter, 67**
- Observability**
 - queue issues, 378
 - spinlocks, 376–377
- Observation tools**
 - kinds, 167–169
 - logging, 133
- Observation tools design, 209**
 - consequences, 212
 - data display, 214–215
 - frequency and length decisions, 210–211
 - histogram buckets case study, 212–214
 - item selection, 209–210
 - overhead, 211–212
 - summary, 215–216
- Observations**
 - matrix multiplication, 52–53
 - network traffic, 94–96

- Offered loads
 - constraints, 11
 - description, 3
 - limits, 4
 - overview, 189–191
 - On-disk buffering, 66–68
 - On-disk databases, 121
 - On-disk read cache, 68
 - On-disk write buffer, 67, 68, 78, 79, 124
 - On the wire, packet, 89, 320–321, 323, 325, 329
 - `open()` system call, 67
 - Oppen, Frederick Burr, 338, 414
 - oprofile** profiling system, 173–176
 - Option** group in software dynamics display, 258–259
 - Order-of-magnitude estimates, 9–10
 - Oscillators for clock, 111–112
 - Overhead
 - IPC trace entries, 232
 - KUtrace tool design, 221
 - observation tools design, 211–212
 - software dynamics display timelines, 264–265
 - trace tools, 197–198, 202, 205, 207
 - Overlapping intervals in aggregate measures, 152–153
 - OUIs (Organizationally Unique Identifiers), 88
 - Ousterhout, John K., 157, 410
 - Out-of-order, 19, 40, 41, 47, 325
- ## P
-
- Packets
 - Ethernet, 87–88
 - overview, 90–91
 - software dynamics display timelines, 263
 - TCP/IP, 90
 - tracing, 234–235
 - Padding field in RPC headers, 99
 - Padegs, Andris, 376, 410
 - `page_fault` routine, 172, 180, 181, 304
 - Page faults
 - in child threads, 293
 - KUtrace tool for, 217, 244
 - mystery25 program, 308
 - page table access, 304
 - `paging_hog` program and, 299, 300–304
 - patches tracing, 236
 - performance counters, 172
 - in `schedtest`, 304
 - Page table, 46, 82, 227, 293, 299, 304
 - Paged memory, 17
 - Paging activity
 - CPU waiting problem, 293
 - memory waiting problem, 299–305
 - paging_hog** program, 299–300, 303
 - Painter, Richard A., 280, 410
 - Parady, Bodo, 409
 - Parallel servers for datacenter transactions, 6–7
 - Parent ID field in RPC headers, 98
 - Parity for memory, 33
 - Partially full intervals in aggregate measures, 153
 - Patches
 - KUtrace tool, 224, 394
 - Linux kernel. *See* Linux kernel patches for KUtrace
 - Patterns in memory access, 31–32
 - Patterns of values in aggregate measures, 151–152
 - Patterson, David A., 15, 408
 - PaX project, Linux, 264, 410
 - Payload, packet, 88, 97, 235, 401
 - PC (program counter), 133
 - profiling tools, 168, 173
 - samples in display timelines, 263
 - SMT, 18
 - trace example, 197–199
 - PCIe (peripheral component interconnect express) bus, 65
 - Pendharkar, Sneha, 407
 - Pentium 4 processor, 18
 - Pentium P5 processor, 20
 - Per-event measurements, aggregating, 150–151
 - Per-function counts and time trace example, 199–202
 - Per-function Gmail trace example, 203–207
 - Per-instance dashboards, 163–164
 - Per-server dashboards, 164
 - Percentiles
 - aggregate measures, 145–147, 150–151
 - long-tail latency, 8–9
 - perf** command, 171–173
 - Performance analysis overview, 269–270
 - Peripheral component interconnect express (PCIe) bus, 65
 - Perl, Sharon E., xxv, 195, 196, 410
 - Physically addressed cache, 37
 - `pid` field for JSON events, 255
 - PID** group in software dynamics display, 259–260
 - PID groups in queue issues, 373
 - Piecewise-linear graphs for aggregate measures, 147, 149
 - ping** method
 - RPC clients, 104
 - RPC logs, 100
 - Pipelines in ACS-1 computer, 17
 - Pipelining instruction, 15, 16, 19, 27, 401
 - PlainSpinLock** class, 366, 374–375
 - Plakal, Manoj, 411
 - Planes in core memory, 32–33
 - Platters in hard disks, 62–63
 - Plausibility check, 52, 165
 - Point events
 - fixed-length trace entries, 392
 - software dynamics display timelines, 262
 - poll** service, 358
 - Ports in TCP/IP, 90

Postprocessing for KUtrace
 details, 249
 eventtospan program, 251–252
 implementation, 223, 225
 JSON format, 254–256
 makeself program, 254
 overview, 249
 rawtoevent program, 250–251
 samptoname_k and samptoname_u programs, 253–254
 spantospan program, 253
 spantotrim program, 253
 summary, 256
 trace entry events, 395

Postprocessing RPC logs, 106

Power-saving mode
 advent of, 20
 changing parameters of, 370
 as idle option, 262, 304
 observation of, 222
 performance disasters and, 161, 162
 states for KUtrace tool, 222
 and timestamps, 21

Power4 chip, 18

Prabhakar, Balaji, 407

Preamble, packet, 20

Precharge cycle, DRAM, 42

Prefetch, cache, 40–41, 57, 61, 401

PreProcessEvent routine, 252

Prescher, Thomas, 409

Primary tasks for queues, 362, 365

/proc pseudofile, 171

process_message routine, 334

Process wakeup in software dynamics display timelines, 262

ProcessEvent routine, 251–252

Readfiles
 logs, 133
 observation tools design, 212
 tools, 167–169

Program counter. See PC

Protocol buffers, 138

Pseudofiles, 171

Pseudorandom RPC IDs, 137

pthread_create routine, 293

Q

Queries
 description, 3
 latency, 4

Queues, waiting for
 complex examples, 370
 CPU frequency, 369–370
 cross-checking, 379–380
 dequeue routine, 365–366
 enqueue routine, 366
 exercises, 380–381
 KUtrace tool, 371–374

latency, 368–369
 load balancing, 377–378
 overview, 361
 primary tasks, 365
 queue depth, 151, 163, 378
 queue structure, 364
 request distribution, 363–364
 root cause, 375–376
 RPC logs, 370–371
 simple examples, 367–368
 spin, 378–379
 spinlocks, 366, 374–377
 summary, 380
 work routine, 367
 worker tasks, 365

queuetest program, 361, 367–368, 370, 375–379

quit method
 RPC clients, 105
 RPC logs, 100

R

Randell, Brian, 410

Random 4KB blocks, reading, 313–314, 316

randomid field for JSON metadata, 255

Ranganatha, Parthasarathy, 408

Raspberry Pi-4B boards, 237

Raw traceblock format, 229–230

rawtoevent program, 250–251

RDTSC (Read timestamp counter instruction), 20, 22–24, 201–202

Read-copy-update (RCU) technique, 353–354

read method
 RPC clients, 104
 RPC logs, 100

Read/write heads in hard disks, 62–64

Readouts in memory, 33

Reads
 disk, 40MB, 310
 disk, caches, 67
 disk, random 4KB blocks, 313–314
 disk, sequential 4K blocks, 311–313
 disk, speed, 68–71
 disk, timing displays, 73–77
 SSD, 40MB, 315
 SSD, timing displays, 80–81

really_get_slow_claim() routine, 244

Receive-side scaling. See RSS

record command, 171–172

recvmmsg service, 358

Refresh, 33, 213, 401

Regions in software dynamics display
 Region 1, controls, 258–259
 Region 2, Y-axis, 259–260
 Region 3, timelines, 260–265
 Region 4, IPC legend, 265
 Region 5, X-axis, 265
 Region 6, save/restore, 265

Register-to-register operations in early CPUs, 16

Regits, William M., 33, 412

ReleaseLock code, 344

Reliable connections, 90

Remote procedure calls. *See* RPCs

report command, 171–172

Requests

arrival times, on master dashboards, 159–160

description, 3

latency times, on master dashboards, 160–161

per-instance dashboards, 163

queue waiting distribution, 363–364

reset method

RPC clients, 104

RPC logs, 100

Resonant frequency, 64

Resources, fundamental, 12, 15, 31, 169, 262, 398

Response times

controlling, 5, 7

CPU business and, 4

exceptionally slow, 95, 134

histogram case study, 212

per-instance dashboards, 163–164

transaction, 401

understanding, 1

wide variety of, 154, 272

Restore, 33, 265, 401

Retransmit, 321, 331, 333, 336

Retransmit timeout, 332, 333, 359

retval field

JSON events, 256

trace entries, 231, 391

Ring buffer, 324, 325, 329

Root cause, 373, 375, 384

Rosenblum, Mendel, 407

Rothenberg, Jeff, 138, 410

Round-robin, 44, 290, 346, 349

Routers, 89

Row access, array, 50, 57

Row access, DRAM, 42

RPC group in software dynamics display, 260

RPC ID, 98

RPC_stats routine, 353–354

rpcid field for JSON events, 255–256

RPCs (Remote procedure calls)

aggregate measures, 150–151

client program, 102–105

client/server systems, 100, 105–106

data, 104, 105, 115

datacenters, 6

disk and network database interaction,

115–125

exercises, 109

headers, 97–99, 104, 105

IDs in logs, 136–137

interference, 11

intermittent performance, 273–277

KUtrace tool design, 222

latencies, 8

logs postprocessing, 106

marker, 97, 98, 104, 105, 235

messages, definition, 96–99

messages, networks, 85–87

messages, server program, 101

messages, TCP/IP, 90

messages, traffic observations, 94–96

mystery21 program, 271–275

network waiting issues. *See* Networks,

waiting for

observations, 107–108

overview, 91–93

packet tracing, 235

queue logs, 370–371

summary, 108–109

trees, 7

RSS, 335

Rubtsov, Artem, 63, 410

Rule of thumb, 69, 76, 211, 212

Ryzen processor chips, 32

S

Saidi, Ali G., 405

Saive, Ravi, 410

Samp button for software dynamics display, 259

Sample servers

connecting, 388–389

hardware, 387–388

samptaname_k program, 253–254, 264

samptaname_u program, 253–254, 264

Sanity checks on dashboards, 164–165

Sankar, Sriram, 407

SATA (Serial AT attachment), 65, 75, 81, 82, 312

Save/Restore region, 265

Scales for aggregate measures, 147

Schedulers

CPU waiting issues, 289–292

Linux kernel patches, 236

Schmidt, John D., 33, 406, 410

Scholz, Hans-Peter, 413

Schwarz, Michael, 409

ScrambledTiming() routine, 42

Scott, Tom, 136, 410

Screw-up, performance

200 usec delay, 327

benchmark dilution, 281

benchmark program measures, 285

code oversight, 353

critical-section locks, 189

Checksum calls, 276

CPU slowdown, 369

deep algorithm changes for, 174

dropped transactions, 377

external interference, 368

- futex inside the critical section, 375, 376
- with gcc compiler, 24, 284
- with gcc optimizer, 24
- hardware glitch, 210
- histogram buckets, 213
- with index updates, 207
- IPC trace entries, 232
- latency under light loads, 161
- from legacy debug code, 214
- lock, delays from, 353
- lock captures, 347
- lock saturation, 345, 348
- misconfigured kernel, 333
- mwaiT instruction timing, 293
- mystery27 program, 341
- no observability, 376, 378
- PlainSpinLock flaw, 375
- profile flaws, 203
- program interference, 368
- queue delay, 373
- retransmit timeout, 332, 333, 359
- scheduler failure, 295, 291, 294, 295
- sequential threads, 270
- server code variability, 272
- slow rdtsc, 201
- slowdown from spinning, 379
- slowdown of Whetstone module, 286
- straight-line code fail, 21
- time between RPCs, 328
- UDP packet bunching, 335
- variable names, misleading by, 347
- variation, misunderstanding, 203
- write data buffering, 124
- writer lock in reader lock, 347
- writing and sync of 40MB, 309
- Search group in software dynamics display, 259**
- Secondary controls in software dynamics display, 265–266**
- Sector, disk. See Disk blocks**
- Security, 17, 87, 91, 225**
- Seek time**
 - disk read speeds, 69, 70, 74
 - experiments involving, 125, 126
 - flash drive delays, 68
 - for 4K block reads, 313, 316
 - paging_hog program, 303
 - plus transfer time, 124
 - SATA bus, 65
 - solid-state drive, 314
- Selective ACK, 332, 333**
- send_to calls, 334**
- SEQ_TCP, 88, 309–315**
- Sequential 4K blocks**
 - reading on disk, 311–313
 - reading on SD, 315
- Serenyi, Denis, 408**
- Serial AT attachment. See SATA**
- server_disk program, 121**
- server4 program**
 - log files, 106, 138
 - network waiting problem, 320–321, 329
 - sample server, 100–101
- Servers**
 - datacenters, 5–7
 - description, 3, 86
 - latency data on master dashboards, 161
 - RPCs, 101
 - sample, 387–389
 - time-of-day clock, 113
- Service level agreements (SLAs), 190–191**
- Service response time commitments, 190–191**
- Services, 3**
- Set access bias, cache, 51**
- Set-associative caches, 36**
- Seurat, Georges, 168**
- Shacham, Nachum, 347, 410**
- Shanbhag, Chandan, 411**
- Shift registers for RPC IDs, 137**
- shortMulX field for JSON metadata, 255**
- shortUnitsX field for JSON metadata, 255**
- show_cpu.html file, 257–258**
- Shrivastav, Vishal, 409**
- Signatures for RPC markers, 97**
- Sigelman, Benjamin H., 7, 94, 189, 195, 196, 410**
- Simultaneous multithreading. See SMT**
- Sine wave, idle delay, 232, 293, 315, 326**
- Sites-Bowen, Connor J., xxv, 71, 414**
- Sites, Richard L., 195, 196, 204, 293, 405, 409, 411**
- Size**
 - cache levels, 43–45
 - cache lines, 38–39
- Skew, execution, 7, 13, 398**
- Skewed distribution, 363–364, 367, 370, 373, 377**
- Skewed requests in queues, 362–364**
- SLAs (service level agreements), 190–191**
- Sleep**
 - AMD Ryzen, 326
 - avoiding, 317
 - Intel C6, 293, 294
 - mwaiT and, 315
 - thread delays from, 297
 - time between RPCs and, 327, 328
 - time to come out of, 313
 - and work completion costs, 379
- Slop**
 - communication, 93, 402
 - disk and network database interaction time, 115–116
 - networks, 93–95
- Slow performance problem**
 - analysis, 286
 - floating-point antagonists, 282–284
 - memory antagonists, 285

- mystery, 280–282
 - overview, 279
 - program, 279–280
 - summary, 286–287
 - Whetstone benchmark, 279–282
 - SMT (simultaneous multithreading), 18**
 - Software**
 - critical sections. *See* Critical sections
 - datacenters, 6–7
 - Software disk access, 66–68**
 - Software dynamics display. *See* Display of software dynamics**
 - Software layers, 5, 6, 86**
 - Software locks**
 - CPU waiting issues, 295–297
 - software dynamics display timelines, 263
 - waiting for. *See* Locks, waiting for
 - Solid-state drives. *See* SSD**
 - Sorting, 73, 138, 162, 199, 215, 252, 254**
 - spantotrim program, 253**
 - spantotrim program, 253**
 - SPEC89 benchmarks, 58, 280**
 - SPEC92 benchmarks, 58**
 - Spector, Alfred Z., 407**
 - Spectre, 225, 409**
 - Speculative execution, 16, 19, 211, 343**
 - Speed concerns overview**
 - datacenter context, 3–5
 - datacenter hardware, 5–6
 - datacenter software, 6–7
 - fundamental resources, 12
 - long-tail latency, 7–9
 - order-of-magnitude estimates, 9–10
 - summary, 12–13
 - thought framework, 9
 - transactions, 11–12
 - Spin with queues, 378–379**
 - Spinlock, fancy, 339, 346, 364, 379**
 - Spinlock, plain, 365, 366, 374–376**
 - Spinlocks**
 - disk and network database interaction, 118
 - networks, 101–102
 - PlainSpinLock Flaw, 374–377
 - queues, 366
 - Spock, S'chn T'gai, 386**
 - spn box in software dynamics display, 266**
 - printf() routine, 137**
 - SRAM (static random-access memory), 33, 34, 402**
 - SSDs (solid-state drives)**
 - banks, 81, 317
 - overview, 64–66
 - writing and sync of 40MB, 314–315
 - Standard deviations in aggregate measures, 145–147**
 - Standard Performance Evaluation Corporation.**
 - See* SPEC
 - Startup**
 - checks at, 164, 165
 - idle CPU time at, 292, 304
 - half-useful principle and, 70, 317
 - log files for, 168
 - perf program, 173
 - recording, 134
 - sort program, 198
 - sync system, 309
 - timestamped log of, 357
 - unusual delays in, 371
 - Starvation, 127, 337, 338, 344, 347, 348**
 - stat command, 171–172**
 - Static random-access memory (SRAM), 33**
 - stats method**
 - RPC clients, 104
 - RPC logs, 100
 - Status field in RPC headers, 99**
 - Stephenson, Pat, 411**
 - Stevens, W. Richard, 321, 407**
 - Stoll, Clifford P., 379, 411**
 - Stoner, M. J., 407**
 - stop command in ktrace_control program, 246**
 - strace command, 176–179**
 - Straight-line code fail measurements, 21**
 - Stress test, 370, 380, 386**
 - Stride, array, 38–40, 42, 43, 46**
 - Strings for logs, 134**
 - Subblock, array, 35, 55, 57–58**
 - Subsystems in datacenters, 6**
 - Superlinear slowdown, 316**
 - Superscalar design**
 - ACS-1 computer, 17
 - microprocessors, 18
 - Supervisory programs in datacenters, 6**
 - Switches and switching**
 - networks, 89
 - traceblocks, 244
 - Sync system call, 309, 310**
 - Synchronization**
 - 40MB on SSD, 314–315
 - disk and network database interaction time, 111–117
 - disk operations, 308–310
 - Synchronous, 6, 128, 403**
 - /sys pseudofiles, 171**
 - Syscall trace entries, 231**
 - System 360/91 design, 17**
 - System under test**
 - defined, 403
 - distortion of, 133, 194, 201, 211, 215
 - framework for examining, 9
 - offered load driving, 128
-
- T**
- T field in trace entries, 231, 391**
 - T1 field in RPC headers, 98**
 - Tacoma Narrows Bridge, 413**
 - Tag, cache, 35–37, 403**
 - Tail latency**
 - description, 4

- frameworks, 9
- histogram time scales, 147–149
- overview, 7–9
- Tails in queue structure, 364
- TCP (Transmission Control Protocol) packet tracing, 234–235
- TCP/IP (Transmission Control Protocol/Internet Protocol)
 - network waiting issues, 321
 - overview, 89–90
- tcpdump** tool
 - network traffic, 187–189
 - network waiting issues, 325–327, 329, 331
 - packet information, 195
 - packet tracing, 235
- Technical term, 5, 124
- Tensor processing units (TPUs), 385
- Terms, glossary of, 397–404
- Tesla, xxv, xxvii
- Test-and-set instructions in lock waiting, 338–339
- Text strings for logs, 134
- Thacker, Charles P., v
- Thoth trace tool, 204
- thousandsX field for JSON metadata, 255
- Thrash, 36, 181, 403
- Threads
 - description, 5
 - locks. *See* Locks, waiting for spinlocks, 101–102
- Thurber, James G., 414
- Time, waiting for, 357
 - inline execution delays, 359
 - periodic work, 357–358
 - summary, 359
 - timeouts, 358
 - timeslicing, 358–359
- Time alignment, 111–117, 118, 155, 266
- Time-base problem in packet tracing, 235
- Time-constrained software
 - deadlines, xxi
 - default delays, 333
 - incentives for, 5
 - KUtrace to observe, 217, 220, 221
 - observation tools, 131, 133, 169, 175, 211
 - vs. offline software, 4
 - with performance issues, xxii
 - timeouts and, 358
 - unacceptable delays for, 327
- time** command, 171
- Time groupings on master dashboards, 161–162
- Time-of-day clock, 112
- Time scales
 - histogram, 147–149
 - trace tools, 196
- Time zones with timestamps, 136
- timealign.cc** program, 115–117
- TimeDiskWrite()** routine, 73
- Timelines**
 - aggregate measures, 143–147
 - software dynamics display, 260–265
- Timeouts**, 358
- Times in data display**, 215
- Timeslicing**, 358–359
- timestamp field for JSON events**, 255
- Timestamps**
 - constant-rate counters, 21
 - KUtrace tool, 219
 - Linux kernel patches, 233
 - logs, 134–136
 - packet tracing, 235
 - per-function Gmail trace example, 205
 - per-instance dashboards, 164
 - queue waiting, 362
 - trace entries, 391
- Timing for memory**, 31–32
- title field for JSON metadata**, 255
- TLBs (translation lookaside buffers)**
 - early CPUs, 17
 - organization, 36–37
- Tools**
 - blktrace, 184–187
 - Dapper, 189–191
 - data to be observed, 169–170
 - design. *See* Observation tools design
 - exercises, 191
 - ftrace, 180–183
 - locktrace, 189
 - ltrace, 179–180
 - mtrace, 183–184
 - observation, 167–169
 - oprofile profiling system, 173–176
 - perf, 171–173
 - /proc and /sys pseudofiles, 171
 - strace, 176–179
 - summary, 191
 - tcpdump, 187–189
 - time, 171
 - top, 170–171
 - Wireshark, 187–189
- top** command, 170–171
- TPUs (tensor processing units)**, 385
- Trace calls**
 - initializing and controlling, 241
 - Insert1, 241–243
 - InsertN, 243–244
- Trace entries**
 - event numbers, 393–395
 - fields, 391
 - fixed-length, 391–392
 - IPC, 232
 - nested, 233
 - syscall, 230–231
 - variable-length, 392–393
- tracebase field for JSON metadata**, 255

Traceblocks

- Linux kernel patches, 229–230
- switching, 244

Traces, 193

- advantages, 193–194
- buffer data structures, 228
- design questions, 194–197
- disadvantages, 194
- initializing and controlling, 241
- KUtrace tool design, 222
- logs, 133
- per-function counts and time example, 199–202
- per-function Gmail example, 203–207
- program counter example, 197–199
- summary, 207–208
- user-mode runtime control, 245

Tracing tools, 167–169

- blktrace, 184–187
- ftrace, 180–183
- locktrace, 189
- ltrace, 179–180
- mtrace, 183–184
- observation tools design, 212
- strace, 176–179

Track, disk, 62–64**Transactions**

- aggregate measures examples, 154–155
- description, 3
- latency, 189–191
- latency variation in disk and network data-base interaction, 128
- speed factors, 11–12

Translation buffers access time, 46**Translation lookaside buffers. See TLBs****Transmission Control Protocol (TCP) packet tracing. See TCP****Transmission Control Protocol/Internet Protocol. See TCP/IP****Transmission delays in networks, 323–327****Transpose method, 55–57****Trees for remote procedure calls, 7****Tufte, Edward R., 214, 411****Turner, Paul, 408****txt box in software dynamics display, 266****Type field in RPC headers, 98**

U

UDP (User Datagram Protocol) packets

- network waiting anomaly, 334–335
- tracing, 234–235

Unaligned references, 36**Unbalanced tasks in queues, 377–378****Underutilization of caches, 46****Unidentified communication time in networks, 93–94****Uniform distribution, 363, 364, 381****Uniform event rates vs. bursty, 142****Uniform requests in queues, 362–364****Units**

- dashboards, 160
- order-of-magnitude, 9

Unpredictable unusual transactions, traces for, 194**Unreasonable offered loads, 11****Update intervals**

- aggregate measures, 152–154
- master dashboards, 160–161

User Datagram Protocol. See UDP**User-facing**

- complex software, 4
- fan outs from, 7
- foreground programs, 5, 157, 158
- latency, understanding, xxiv
- live load, 11, 194
- performance disasters, 68
- RPCs, 131
- unacceptable delays, 122, 125

User-mode code

- RPCs, 92, 320
- trace entry events, 395

User-mode library for KUtrace tool, 219, 222**User-mode runtime control for KUtrace, 245**

- control interface to loadable modules, 247
- kutrace_control program, 246
- kutrace_lib library, 246–247
- summary, 247
- tracing control, 245

UTC (coordinated universal time) for timestamps, 136

V

Vahdat, Amin M., 407, 408**Valgrind, 189, 411****Vampire taps, 87****Variable-length trace entries, 392–393****Variance, xxii, 1, 168, 323****Variation**

- CPU frequency, 369
- disk blocks and, 75
- execution skew as, 7
- finding sources of, 1
- transaction latency, 95, 111, 128
- mystery2.cc, 47, 346, 349
- 99th percentile, 151
- profiling and, 203
- RPC logging showing, 272
- run-to-run, 52
- 3X, 42

VAX-11/780, 407**version field for JSON metadata, 255****VLAN packets, 90–91****Virtual machine, 220, 385****Virtual memory, 15, 17, 36, 46, 82, 170, 306**

Virtually addressed cache, 37
vmalloc routine, 227
 Volatile variables in CPU measurements, 25–26, 29

W

wait command in `kutrace_control` program, 246
waitMsec routine, 328
 Wall, David W., 406
 Walpole, Jonathan, 409
 Wang, Han, 409
 Watchdog timers, 358
 Wear-leveling in SSDs, 66
 Weatherspoon, Hakim, 409
 Weaveworks, 411
 Whetstone benchmark, 25, 173, 191, 279, 281, 286
 Wikimedia 2005, 63, 414
 Wikimedia 2006, 65, 414
 Wikimedia 2008, 143, 411
 Wikimedia 2010, 32, 413
 Wikimedia 2012, 62, 414
 Wikimedia 2013, 62, 414
 Wikimedia 2016, 33, 413
 Wikimedia 2020a, 88, 414
 Wikimedia 2020b, 338, 414
 Wikimedia 2021, 111, 414
 Wikipedia 2019a, 17, 411
 Wikipedia 2020a, 16, 411
 Wikipedia 2020b, 17, 411
 Wikipedia 2020c, 17, 411
 Wikipedia 2020d, 18, 411
 Wikipedia 2020e, 18, 412
 Wikipedia 2020f, 33, 412
 Wikipedia 2020g, 290, 412
 Wikipedia 2020h, 345, 412
 Wikipedia 2021a, 16, 412
 Wikipedia 2021b, 16, 412
 Wikipedia 2021c, 17, 412
 Wikipedia 2021d, 17, 412
 Wikipedia 2021e, 17, 412
 Wikipedia 2021f, 18, 412
 Wikipedia 2021g, 18, 412
 Wikipedia 2021h, 18, 412
 Wikipedia 2021i, 18, 412
 Wikipedia 2021j, 18, 412
 Wikipedia 2021k, 20, 412
 Wikipedia 2021l, 18, 412
 Wikipedia 2021m, 33, 412

Wikipedia 2021n, 87, 412
 Wikipedia 2021o, 112, 412
 Wikipedia 2021p, 136, 413
 Wikipedia 2021q, 137, 413
 Wikipedia 2021r, 138, 413
 Wikipedia 2021s, 143, 413
 Wikipedia 2021t, 186, 413
 Wikipedia 2021u, 199, 413
 Wikipedia 2021v, 320, 413
 Wikipedia 2021w, 320, 413
 Wikipedia 2021x, 345, 413
 Wikipedia 2021y, 347, 413
 Wikipedia 2021z, 354, 413
 Williams, Don, 190, 413
 Wireshark, 187–189, 235, 319
 Words

- instruction sets, 16
- memory, 32–33

Worker tasks in queues, 362, 365
worker_thread processes, 341–344

Write cycle, 20, 65, 66

write method
 RPC clients, 104
 RPC logs, 100

Writes
 disk, 40MB, 308–310
 disk, buffering, 67
 disk, speed, 71–73
 disk, timing display, 77–80
 SSD, 40MB, 314–315
 SSD, cycles, 65
 SSD, timing display, 82

X

X-axis in software dynamics display, 265
 Xeon processor, 18
 Xerox PARC, 87

Y

Y-axis in software dynamics display, 259–260
 Yarom, Yuval, 409
Ychars box in software dynamics display, 266
 Yin, Zi, 407
 YouTube 2016, 64, 413

Z

Zwaenepoel, Wily, xxv
 Zwergelstern, 414