

Wydajność i optymalizacja kodu

Istota dynamiki działania
oprogramowania

Richard L. Sites



Przedmowa Luiz André Barroso, Google Fellow

Tytuł oryginału: Understanding Software Dynamics

Tłumaczenie: Tomasz Walczak

ISBN: 978-83-283-9515-2

Authorized translation from the English language edition, entitled Understanding Software Dynamics, 1st Edition by Richard Sites, published by Pearson Education, Inc, publishing as Addison-Wesley Professional, Copyright © 2022 Pearson Education, Inc.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from Pearson Education, Inc.

Polish language edition published by Helion S.A., Copyright © 2023.

Wszelkie prawa zastrzeżone. Nieautoryzowane rozpowszechnianie całości lub fragmentu niniejszej publikacji w jakiegokolwiek postaci jest zabronione. Wykonywanie kopii metodą kserograficzną, fotograficzną, a także kopiowanie książki na nośniku filmowym, magnetycznym lub innym powoduje naruszenie praw autorskich niniejszej publikacji.

Wszystkie znaki występujące w tekście są zastrzeżonymi znakami firmowymi bądź towarowymi ich właścicieli.

Autor oraz wydawca dołożyli wszelkich starań, by zawarte w tej książce informacje były kompletne i rzetelne. Nie biorą jednak żadnej odpowiedzialności ani za ich wykorzystanie, ani za związane z tym ewentualne naruszenie praw patentowych lub autorskich. Autor oraz wydawca nie ponoszą również żadnej odpowiedzialności za ewentualne szkody wynikłe z wykorzystania informacji zawartych w książce.

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<https://helion.pl/user/opinie/roprwy>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 230 98 63

e-mail: helion@helion.pl

WWW: <https://helion.pl> (księgarnia internetowa, katalog książek)

Printed in Poland.

- Kup książkę
- Poleć książkę
- Oceń książkę

- Księgarnia internetowa
- **Lubię to!** » Nasza społeczność

Spis treści

Przedmowa	15
Wprowadzenie	16
Podziękowania	21
O autorze	22
Część I Pomiary	23
1. Mój program działa zbyt wolno	25
1.1. Kontekst centrum danych	25
1.2. Sprzęt w centrach danych	27
1.3. Oprogramowanie w centrum danych	28
1.4. Latencja z długiego ogona rozkładu	30
1.5. Model myślenia	32
1.6. Szacowanie rzędu wielkości	32
1.7. Dlaczego transakcje działają powoli?	33
1.8. Pięć podstawowych zasobów	35
1.9. Podsumowanie	35
2. Pomiary procesorów	37
2.1. Trochę historii	38
2.2. Obecna sytuacja	41
2.3. Pomiar latencji instrukcji add	42
2.4. Niepowodzenie z prostym, sekwencyjnym kodem	44
2.5. Niepowodzenia z prostą pętlą, kosztami wykonywania pętli i kompilatorem optymalizującym	44
2.6. Niepowodzenie z martwą zmienną	47
2.7. Lepsza pętla	48
2.8. Zmienne zależne	49
2.9. Faktyczna latencja wykonywania	49
2.10. Więcej niuansów	50
2.11. Podsumowanie	51
Ćwiczenia	51
3. Pomiar pamięci	53
3.1. Pomiar czasu dostępu do pamięci	53
3.2. Pamięć	54
3.3. Struktura pamięci podręcznej	56

3.4. Wyrównanie danych	59
3.5. Struktura bufora TLB	60
3.6. Pomiary	61
3.7. Pomiar wielkości wiersza pamięci podręcznej	62
3.8. Problem: wstępne wczytywanie wiersza $N + 1$	64
3.9. Odczyt uzależniony od poprzedniej operacji	65
3.10. Nielosowy dostęp do pamięci DRAM	66
3.11. Pomiar łącznej wielkości każdego poziomu pamięci podręcznej	68
3.12. Pomiar stopnia wielodrożności pamięci podręcznej na poszczególnych poziomach	70
3.13. Czas dostępu do bufora TLB	71
3.14. Niepełne wykorzystanie pamięci podręcznej	71
3.15. Podsumowanie	71
Ćwiczenia	72
4. Interakcje procesora i pamięci	74
4.1. Interakcje związane z pamięcią podręczną	75
4.2. Dynamika prostego mnożenia macierzy	76
4.3. Szacunki	77
4.4. Inicjowanie, kontrola wyników i obserwacja	78
4.5. Początkowe wyniki	78
4.6. Szybsze mnożenie macierzy metodą transpozycji	81
4.7. Szybsze mnożenie macierzy z wykorzystaniem podbloków	83
4.8. Obliczenia z uwzględnianiem pamięci podręcznej	84
4.9. Podsumowanie	85
Ćwiczenia	85
5. Pomiar dysków twardych i nośników SSD	86
5.1. Dyski twarde	86
5.2. Nośniki SSD	90
5.3. Dostęp do dysku w programie i buforowanie danych na dysku	92
5.4. Jak szybki jest odczyt danych z dysku?	94
5.5. Nieco prostych obliczeń	97
5.6. Jak szybki jest zapis danych na dysku?	99
5.7. Wyniki	99
5.8. Odczyt danych z dysku	100
5.9. Zapis danych na dysku	104
5.10. Odczyt danych z nośnika SSD	107
5.11. Zapis danych na nośniku SSD	109

5.12. Wiele transferów	109
5.13. Podsumowanie	110
Ćwiczenia	111
6. Pomiar dotyczące sieci	113
6.1. Ethernet	115
6.2. Koncentratory, przełączniki i routery	117
6.3. Protokół TCP/IP	118
6.4. Pakiety	119
6.5. Wywołania RPC	120
6.6. Niezidentyfikowany czas	122
6.7. Obserwowanie ruchu w sieci	123
6.8. Definicja przykładowego komunikatu RPC	126
6.9. Projekt rejestrowania zdarzeń	129
6.10. Przykładowy system klient-serwer oparty na wywołaniach RPC	130
6.11. Przykładowy program serwera	131
6.12. Blokady wirujące	132
6.13. Przykładowy program klienta	133
6.14. Pomiar jednego przykładowego wywołania RPC między klientem a serwerem	136
6.15. Przetwarzanie końcowe dzienników wywołań RPC	137
6.16. Obserwacje	137
6.17. Podsumowanie	139
Ćwiczenia	140
7. Dyskowa baza danych i jej interakcje z siecią	142
7.1. Wyrównywanie pomiarów czasu	142
7.2. Wiele klientów	149
7.3. Blokady wirujące	149
7.4. Pierwszy eksperyment	150
7.5. Baza danych na dysku	153
7.6. Drugi eksperyment	153
7.7. Trzeci eksperyment	158
7.8. Rejestrowanie informacji	160
7.9. Wyjaśnienie zmienności latencji transakcji	160
7.10. Podsumowanie	161
Ćwiczenia	162

Część II Obserwacja	163
8. Rejestrowanie zdarzeń	165
8.1. Narzędzia do obserwacji	165
8.2. Rejestrowanie zdarzeń	166
8.3. Podstawowe mechanizmy rejestrowania zdarzeń	166
8.4. Rozbudowane rejestrowanie zdarzeń	167
8.5. Znaczniki czasu	168
8.6. Identyfikatory wywołań RPC	169
8.7. Formaty plików dziennika	170
8.8. Zarządzanie plikami dziennika	171
8.9. Podsumowanie	172
9. Miary zagregowane	173
9.1. Zdarzenia występujące jednostajnie i seryjnie	174
9.2. Mierzone okresy	175
9.3. Oś czasu	175
9.4. Dalsze podsumowywanie osi czasu	177
9.5. Skale czasowe dla histogramów	179
9.6. Agregowanie pomiarów dotyczących zdarzeń	182
9.7. Zmiany wzorców wartości w czasie	183
9.8. Czas między aktualizacjami	184
9.9. Przykładowe transakcje	186
9.10. Podsumowanie	187
10. Panele kontrolne	189
10.1. Przykładowa usługa	189
10.2. Przykładowe panele kontrolne	191
10.3. Główny panel kontrolny	192
10.4. Panele kontrolne instancji	196
10.5. Panele kontrolne serwerów	197
10.6. Testy poprawności	197
10.7. Podsumowanie	198
Ćwiczenia	199
11. Inne dostępne narzędzia	200
11.1. Rodzaje narzędzi do obserwacji	200
11.2. Obserwowane dane	202
11.3. Polecenie top	204
11.4. Pseudopliki /proc i /sys	204
11.5. Polecenie time	205

11.6.	Polecenie perf	205
11.7.	Narzędzie oprofile, profiler procesora	207
11.8.	Narzędzie strace, wywołania systemowe	211
11.9.	Narzędzie ltrace, wywołania bibliotek języka C w procesorze	214
11.10.	Narzędzie ftrace, śledzenie funkcji jądra w procesorze	214
11.11.	Operacje malloc i free, narzędzie mtrace	217
11.12.	Śledzenie operacji dyskowych, narzędzie blktrace	219
11.13.	Śledzenie sieci, tcpdump i Wireshark	222
11.14.	Blokady sekcji krytycznych, narzędzie locktrace	223
11.15.	Oferowane obciążenie, wywołania wychodzące i latencja transakcji	224
11.16.	Podsumowanie	225
	Ćwiczenia	226
12.	Ślady	227
12.1.	Zalety śledzenia	227
12.2.	Wady śledzenia	228
12.3.	Trzy pytania na start	229
12.4.	Przykład: jeden z pierwszych śladów licznika programu	232
12.5.	Przykład: liczba instrukcji i czas na funkcję	234
12.6.	Studium przypadku: ślady poszczególnych funkcji w serwisie Gmail	238
12.7.	Podsumowanie	243
13.	Zasady projektowania narzędzi do obserwacji	244
13.1.	Co obserwować?	244
13.2.	Jak często i jak długo?	245
13.3.	Jakie koszty są dopuszczalne?	246
13.4.	Konsekwencje projektowe	247
13.5.	Studium przypadku: kubelki w histogramie	247
13.6.	Projektowanie sposobu wyświetlania danych	249
13.7.	Podsumowanie	251
Część III	Narzędzie KUtrace	253
14.	KUtrace: cele, projekt, implementacja	255
14.1.	Ogólne informacje	255
14.2.	Cele	256
14.3.	Projekt	257
14.4.	Implementacja	259
14.5.	Patche i moduł jądra	260
14.6.	Program sterujący	261
14.7.	Przetwarzanie końcowe	261

14.8.	Uwagi na temat bezpieczeństwa	261
14.9.	Podsumowanie	262
15.	KUtrace: patche jądra Linuksa	263
15.1.	Struktury danych bufora śladu	264
15.2.	Format surowych bloków śladu	264
15.3.	Rekordy śladu	267
15.4.	Rekordy z liczbą instrukcji na cykl (I/C)	268
15.5.	Znaczniki czasu	269
15.6.	Numery zdarzeń	270
15.7.	Zagnieżdżone rekordy śladu	270
15.8.	Kod	270
15.9.	Śledzenie pakietów	271
15.10.	Patche dla procesorów x86-64 firm AMD i Intel	273
15.11.	Podsumowanie	274
	Ćwiczenia	275
16.	KUtrace: wczytywany moduł dla systemu Linux	276
16.1.	Struktury danych interfejsu jądra	276
16.2.	Wczytywanie i zwalnianie modułu	277
16.3.	Inicjowanie śledzenia i sterowanie nim	278
16.4.	Implementacja wywołań do generowania śladu	278
16.5.	Insert1	278
16.6.	InsertN	281
16.7.	Przełączanie się do nowego bloku	281
16.8.	Podsumowanie	282
17.	KUtrace: sterowanie w trybie użytkownika	283
17.1.	Sterowanie procesem śledzenia	283
17.2.	Samodzielny program kutrace_control	284
17.3.	Podstawowa biblioteka kutrace_lib	285
17.4.	Interfejs do sterowania wczytywanym modułem	285
17.5.	Podsumowanie	286
18.	Przetwarzanie końcowe w narzędziu KUtrace	287
18.1.	Szczegółowe omówienie przetwarzania końcowego	287
18.2.	Program rawtoevent	288
18.3.	Program eventtospan	290
18.4.	Program spantotrim	292
18.5.	Program spantospan	292
18.6.	Programy samptoname_k i samptoname_u	293

18.7.	Program makeself	293
18.8.	Format plików JSON w narzędziu KUtrace	293
18.9.	Podsumowanie	296
19.	KUtrace: wyświetlanie dynamiki działania oprogramowania	297
19.1.	Wprowadzenie	297
19.2.	Obszar 1 — kontrolki	298
19.3.	Obszar 2 — oś y	300
19.4.	Obszar 3 — oś czasu	301
19.5.	Obszar 4 — legenda dotycząca liczby instrukcji na cykl	307
19.6.	Obszar 5 — oś x	307
19.7.	Obszar 6 — zapisywanie i wczytywanie	307
19.8.	Kontrolki pomocnicze	308
19.9.	Podsumowanie	309
Część IV	Wnioskowanie	311
20.	Na co zwracać uwagę?	313
20.1.	Wprowadzenie	313
21.	Wykonywanie za dużej ilości kodu	316
21.1.	Wprowadzenie	316
21.2.	Program	317
21.3.	Zagadka	317
21.4.	Pomiary i wnioski	318
21.5.	Rozwiązanie zagadki	322
21.6.	Podsumowanie	323
22.	Powolne wykonywanie kodu	324
22.1.	Wprowadzenie	324
22.2.	Program	325
22.3.	Zagadka	325
22.4.	Konkurujący program z operacjami zmiennoprzecinkowymi	328
22.5.	Konkurujący program korzystający z pamięci	330
22.6.	Rozwiązanie zagadki	331
22.7.	Podsumowanie	332
23.	Oczekiwanie na procesor	334
23.1.	Program	334
23.2.	Zagadka	335
23.3.	Pomiary i wnioski	335
23.4.	Zagadka numer 2	336
23.5.	Rozwiązanie zagadki numer 2	338

23.6. Dodatkowa zagadka	341
23.7. Podsumowanie	343
Ćwiczenia	343
24. Oczekiwanie na pamięć	344
24.1. Program	344
24.2. Zagadka	345
24.3. Pomiary i wnioski	345
24.4. Zagadka numer 2 — dostęp do tablicy stron	349
24.5. Wyjaśnienie zagadki numer 2	350
24.6. Podsumowanie	351
Ćwiczenia	351
25. Oczekiwanie na dysk	352
25.1. Program	352
25.2. Zagadka	353
25.3. Pomiary i wnioski	353
25.4. Odczyt 40 MB	356
25.5. Odczyt sekwencyjnych bloków po 4 KB	357
25.6. Odczyt losowych bloków po 4 KB	359
25.7. Zapis i synchronizacja 40 MB na nośniku SSD	361
25.8. Odczyt 40 MB z nośnika SSD	361
25.9. Dwa programy jednocześnie używające dwóch plików	362
25.10. Wyjaśnienie zagadek	364
25.11. Podsumowanie	364
Ćwiczenia	365
26. Oczekiwanie na sieć	366
26.1. Wprowadzenie	367
26.2. Programy	368
26.3. Eksperyment numer 1	368
26.4. Zagadka z eksperymentu numer 1	369
26.5. Pomiary i wnioski z eksperymentu numer 1	371
26.6. Eksperyment numer 1. A co z czasem między wywołaniami RPC?	375
26.7. Eksperyment numer 2	377
26.8. Eksperyment numer 3	377
26.9. Eksperyment numer 4	378
26.10. Wyjaśnienie zagadek	381
26.11. Dodatkowa anomalia	382
26.12. Podsumowanie	384

27. Oczekiwanie na blokady	385
27.1. Wprowadzenie	385
27.2. Program	390
27.3. Eksperyment numer 1 — długi czas utrzymywania blokady	393
27.3.1. Proste zajmowanie blokady	394
27.3.2. Nasycenie blokady	394
27.4. Zagadki w eksperymencie numer 1	395
27.5. Pomiary i wnioski w eksperymencie numer 1	395
27.5.1. Zawłaszczenie blokady	397
27.5.2. Zagłodzenie w oczekiwaniu na blokadę	397
27.6. Eksperyment numer 2 — rozwiązanie problemu zawłaszczania blokady	398
27.7. Eksperyment numer 3 — rozwiązanie problemu rywalizacji przez zastosowanie wielu blokad	399
27.8. Eksperyment numer 4 — rozwiązanie problemu rywalizacji o blokadę dzięki mniejszej ilości pracy przy zajętej blokadzie	401
27.9. Eksperyment numer 5 — eliminowanie rywalizacji o blokadę dzięki zastosowaniu techniki RCU dla panelu kontrolnego	402
27.10. Podsumowanie	404
28. Oczekiwanie na podstawie czasu	406
28.1. Okresowe wykonywanie pracy	406
28.2. Limity czasu	407
28.3. Podział czasu	408
28.4. Wewnętrzne opóźnienia w wykonywaniu	408
28.5. Podsumowanie	409
29. Oczekiwanie na kolejki	410
29.1. Wprowadzenie	410
29.2. Rozkład żądań	412
29.3. Struktura kolejki	413
29.4. Zadania robocze	414
29.5. Zadanie główne	414
29.6. Operacje Dequeue	415
29.7. Operacja Enqueue	415
29.8. Klasa blokady wirującej	415
29.9. Procedura odpowiedzialna za „pracę”	416
29.10. Proste przykłady	416
29.11. Co mogło pójść nie tak?	418
29.12. Częstotliwość procesora	418

29.13. Złożone przykłady	420
29.14. Oczekiwanie na procesory — dziennik wywołań RPC	420
29.15. Analiza oczekiwania na procesor za pomocą narzędzia KUTrace	421
29.16. Błąd w klasie PlainSpinLock	424
29.17. Źródłowa przyczyna	426
29.18. Poprawiona klasa PlainSpinLock zapewniająca obserwowalność	427
29.19. Równoważenie obciążenia	427
29.20. Zapewnianie obserwowalności długości kolejki	428
29.21. Aktywne oczekiwanie na końcu	429
29.22. Jeszcze jedna usterka	430
29.23. Dokładne sprawdzanie	430
29.24. Podsumowanie	431
Ćwiczenia	431
30. Podsumowanie	433
30.1. Czego udało Ci się nauczyć?	433
30.2. Czego nie omówiłem?	435
30.3. Dalsze kroki	436
30.4. Podsumowanie (całej książki)	436
A. Przykładowe serwery	439
A.1. Sprzęt z przykładowych serwerów	439
A.2. Łącza serwerów	441
B. Rekordy śladu	442
B.1. Rekordy śladu o stałej długości	443
B.2. Rekordy o zmiennej długości	443
B.3. Numery zdarzeń	444
B.3.1. Zdarzenia wstawiane przez patche narzędzia KUTrace dla jądra	445
B.3.2. Zdarzenia wstawiane przez kod trybu użytkownika	446
B.3.3. Zdarzenia wstawiane przez kod przetwarzania końcowego	447
Literatura	448
Słowniczek	456

Rozdział 1.

Mój program działa zbyt wolno

Ktoś wchodzi do mojego gabinetu i mówi: „Mój program działa zbyt wolno”. Po chwili przerwy pytam: „A jak wolno powinien działać?”

Dobry programista ma gotową odpowiedź na takie pytanie, ponieważ potrafi opisać zadanie do wykonania i oszacować, ile czasu powinna zająć każda jego część. Taka osoba może odpowiedzieć tak: „To zapytanie do bazy danych sprawdza 10 000 rekordów, z których mniej więcej 1000 okazuje się potrzebnych. Każdy dostęp powinien zajmować około 10 milisekund, a używanych jest 20 dysków, tak więc 10 000 operacji dostępu powinno w sumie zająć mniej więcej 5 sekund. Sieć nie jest obciążona, a procesor i pamięć są używane w niewielkich i prostych operacjach znacznie szybszych niż dostęp do dysku. Całe zapytanie zajmuje około 15 sekund, co jest zbyt długim czasem”.

Mniej staranny programista może odpowiedzieć tak: „Pracowałem całą noc i napisałem 1000 wierszy kodu z użyciem wielu istniejących bibliotek. Kod działa, ale zapytanie zajmuje 15 sekund, a chciałem uzyskać poziom $1/10$ sekundy. Jedna z bibliotek musi być zbyt wolna. Jak mogę ją wykryć?”. Taki programista nie ma pojęcia, czy zasadne jest oczekiwanie czasu na poziomie $1/10$ sekundy. Nie wie też, ile czasu powinno zajmować każde wywołanie kodu z biblioteki lub czy poprawnie korzysta z bibliotek. Nie ma też zaprojektowanego mechanizmu obserwowania dynamiki działania kodu, aby ustalić, na co przeznaczony jest czas programu. W tej książce omawiam wszystkie te zagadnienia.

1.1. Kontekst centrum danych

W tym miejscu wprowadzam pojęcia i koncepcje ze złożonego środowiska oprogramowania. Twoje środowisko może być znacznie prostsze, ale opisywane kwestie można prawie dokładnie do niego przenieść. Używana tu terminologia pochodzi ze świata centrów danych, jednak te same zagadnienia występują również w kontekście baz danych, komputerów stacjonarnych, pojazdów, gier, a także w innych środowiskach z ograniczeniami czasowymi.

Transakcja, zapytanie lub *żądanie* to komunikat wejściowy dla systemu komputerowego przetwarzany jako odrębna jednostka pracy. Każdy komputer przetwarzający transakcje nazywam *serwerem*. *Latencja* lub *czas odpowiedzi* dla transakcji to czas między jej wysłaniem a otrzymaniem wyniku. *Oferowane obciążenie* to liczba transakcji przesyłanych na sekundę. Jeśli ta wartość przekracza liczbę transakcji przetwarzanych na sekundę, następuje wydłużenie czasu reakcji (czasem bardzo znaczne). *Usługa* to kolekcja programów obsługujących transakcje jednego rodzaju. Duże centra danych przetwarzają jednocześnie transakcje dziesiątek różnych usług, a każda usługa ma inne oferowane obciążenie i inne cele związane z latencją.

Latencja transakcji nie jest stała. Ma rozkład prawdopodobieństwa oparty na tysiącach transakcji wykonywanych na sekundę. *Latencje z długiego ogona* dotyczą najwolniejszych transakcji z rozkładu. Prosty sposób na zestawienie takich transakcji jest ustalenie latencji dla percentyla 99%, czyli czasu przekraczanego przez najwolniejsze 1% wszystkich transakcji. Jeśli oferowane obciążenie wynosi 5000 transakcji na sekundę, ten 1% to 50 transakcji.

Dynamika działania programu lub ich grupy oznacza ich aktywność w czasie — jakie fragmenty kodu są wykonywane w różnych momentach, na co oczekują, ile pamięci zajmują i jak poszczególne programy wpływają na siebie nawzajem. Programiści wyobrażają sobie prostą dynamikę działania programu, jednak w rzeczywistości program może, od czasu do czasu, zachowywać się zupełnie inaczej i pracować znacznie wolniej, niż tego oczekujemy. Jeśli uda się zaobserwować prawdziwą dynamikę, można będzie dostosować obraz w umyśle. Zwykle pozwala to na poprawę wydajności kodu za pomocą prostych zmian.

Przedmiotem zainteresowania są tu transakcje z udziałem użytkownika realizowane w złożonym oprogramowaniu, na przykład w centrum danych do obsługi telefonów komórkowych. Ważne są transakcje, które zwykle są realizowane szybko, ale czasami wymagają znacznie więcej czasu, na tyle, że użytkownik końcowy odczuwa irytujące opóźnienie. W centrach danych budżet na sprzęt dla każdej usługi jest często zależny od tego, ile transakcji na sekundę każdy serwer może obsługiwać. Ta docelowa liczba jest określana *empirycznie*, przez zwiększanie oferowanego obciążenia do momentu przekroczenia jakiegoś progowego ograniczenia czasowego dla latencji z długiego ogona w rozkładzie. Następnie ustalone jest nieco niższe obciążenie docelowe.

Jeśli uda Ci się zrozumieć sytuację, a następnie zmniejszyć liczbę zbyt długich transakcji, ten sam sprzęt będzie mógł obsłużyć wyższe obciążenie i nadal realizować cele związane z latencją z długiego ogona. Nie będzie to wymagało ponoszenia dodatkowych kosztów. Takie rozwiązanie zapewnia duże korzyści finansowe. Uzdolniony i mający trochę szczęścia inżynier wydajności może czasem wprowadzić w oprogramowaniu prostą zmianę, która da oszczędności warte jego 10-letniego wynagrodzenia. Firmy i klienci cenią takich ludzi.

Oprogramowanie, w którym transakcje mają ograniczenia czasowe, znacznie różni się od oprogramowania wsadowego lub działającego w trybie offline (a także od większości oprogramowania testowego). Ważną miarą w oprogramowaniu z transakcjami jest czas odpowiedzi, natomiast w oprogramowaniu wsadowym istotną miarą jest zazwyczaj wydajność wykorzystania sprzętu. W transakcjach ważny jest nie tyle *średni* czas odpowiedzi, co najdłuższe czasy, czyli latencja z długiego ogona.

W centrach danych preferowana jest zwykle wyższa średnia latencja z niższą latencją z długiego ogona niż odwrotna sytuacja. Większość osób dojeżdżających do pracy ma podobne preferencje i woli o kilka minut dłuższą drogę, która zawsze zajmuje tyle samo czasu, niż krótszą, na której występują godzinne opóźnienia.

W oprogramowaniu wsadowym średnie obciążenie procesora na poziomie 98% może być korzystne. W oprogramowaniu transakcyjnym 98% obciążenia to *katastrofa*, a nawet średnie obciążenie na poziomie 50% może być zbyt wysokim poziomem, ponieważ skutkuje długimi czasami odpowiedzi, gdy oferowane obciążenie rośnie skokowo na kilka sekund do poziomu trzykrotnie przekraczającego średnie. Gdy w 2004 roku trafiłem do firmy Google, procesory w centrum danych średnio były zajęte przez 9% i bezczynne przez 91% czasu. Obciążenie na poziomie 9% było zbyt niskie. Zwiększenie go do 18% bez wydłużania latencji z długiego ogona pozwoliło podwoić efektywność wszystkich centrów danych. Ponowne podwojenie poziomu, do 36%, byłoby korzystne, jednak kolejne podwojenie, do 72%, prawdopodobnie naruszyłyby zbyt wiele ograniczeń czasowych transakcji.

W trakcie analizowania wydajności złożonego oprogramowania transakcyjnego zakładałem w tej książce, że badane programy działają poprawnie i ich średnia szybkość pracy jest wystarczająco wysoka. Nie omawiam tu projektowania lub debugowania oprogramowania. Nie staram się też zrozumieć lub zwiększyć średniej wydajności. Ponadto przyjmuję, że transakcje, które zawsze działają powoli, zostały zidentyfikowane i poprawione w środowiskach testowych lub diagnostycznych w trybie offline, gdzie nie ma ograniczeń czasowych. Dlatego do zbadania pozostały już tylko okazjonalnie powolne transakcje. Skupiam się na mechanizmach, które powodują długi czas pracy takich transakcji, na tym, jak zaobserwować te mechanizmy, a także na tym, jak zinterpretować poczynione obserwacje.

Gdy używasz telefonu komórkowego do wysyłania wiadomości tekstowych, czytania wpisów, przeszukiwania internetu, przeglądania mapy, strumieniowania filmu, używania aplikacji, a nawet wybierania numeru telefonu, gdzieś działa centrum danych, które reaguje na Twoje żądania. Jeśli odpowiedzi są irytująco powolne, a jakaś konkurencyjna aplikacja lub usługa działa szybciej, może to spowodować, że przejdziesz na nowe rozwiązanie, a przynajmniej zaczniesz rzadziej korzystać z powolnego produktu. W ekosystemie rozwiązań z ograniczeniami czasowymi każdy ma motywację (często finansową) do tego, by ograniczyć irytujące opóźnienia. Jednak niewiele osób potrafi to zrobić.

Ta książka ma sprawić, że takich osób będzie więcej.

1.2. Sprzęt w centrach danych

W dużych centrach danych w budynku znajduje się mniej więcej 10 000 serwerów. Każdy taki serwer to komputer o wielkości zbliżonej do komputera stacjonarnego, ale bez obudowy. Około 50 płyt serwerów jest zamontowanych w szafce, a mniej więcej 200 szafek znajduje się w bardzo dużym pomieszczeniu. Typowy serwer ma od 1 do 4 gniazd na procesory, z których każdy ma po 4 – 50 rdzeni, masę¹ pamięci RAM, kilka dysków twardych lub nośników SSD

¹ To pojęcie z dziedziny informatyki oznaczające 10^{12} .

i połączenie z siecią centrum danych działającą w topologii switching fabric, co umożliwia każdemu serwerowi komunikację ze wszystkimi pozostałymi. Ponadto przynajmniej niektóre serwery mogą komunikować się z internetem, a tym samym z Twoim telefonem. Poza budynkiem znajdują się duże generatory, które zapewniają energię całemu budynkowi (w tym klimatyzacji) na kilka dni lub tygodni na wypadek awarii zasilania. Wewnątrz umieszczone są akumulatory, które potrafią podtrzymać pracę serwerów i przełączników sieciowych przez kilkadziesiąt sekund do czasu uruchomienia generatorów.

Na każdym serwerze działa wiele programów. Zwykle nie ma uzasadnienia biznesowego przydzielanie niektórych serwerów tylko do obsługi poczty elektronicznej, innych tylko do kafelek map, a jeszcze innych do obsługi komunikatorów. Zamiast tego na każdym serwerze pracuje wiele programów, a każdy z nich zwykle ma wiele *wątków*. Na przykład program serwera poczty elektronicznej może mieć 100 wątków roboczych przetwarzających jednocześnie żądania od tysięcy użytkowników, z których większość pisze lub czyta wiadomości. Wiele tych wątków może być aktywnych i oczekiwać na dostęp do dysku lub inne warstwy oprogramowania. Wątki robocze przyjmują przychodzące żądania, realizują je, odpowiadają, a następnie przechodzą do następnego oczekującego żądania od innego użytkownika. W godzinach z największym obciążeniem prawie wszystkie wątki robocze są zajęte. W godzinach, gdy aktywność użytkowników jest mniejsza, przynajmniej połowa wątków jest bezczynna i oczekuje na zadania. W prawie wszystkich skalach czasowych (mikrosekund, milisekund, sekund, minut) występuje ciągły cykl wzrostów i spadków poziomu oferowanego obciążenia. Istnieje nawet siedmiodniowy cykl z niższą aktywnością w soboty i niedziele (zgodnie z dniami roboczymi w zachodnim świecie).

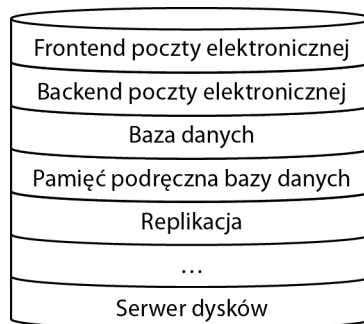
Aby móc kontrolować czas odpowiedzi, ważne jest mieć wolne zasoby sprzętowe dostępne dla transakcji z udziałem użytkowników, ponieważ obciążenie ze strony użytkowników okresowo skokowo rośnie z powodu wydarzeń w świecie fizycznym. Jednak ekonomiczne jest też wykonywanie niekomunikujących się z użytkownikami programów wsadowych, uruchamianych, gdy procesory są bezczynne. Obok komunikujących się z użytkownikami programów pierwszego planu i uruchamianych w tle programów wsadowych na każdym serwerze zawsze działa kilka programów nadzorujących, które śledzą, jak obciążony jest serwer, ile błędów na nim występuje, ile wolnej pamięci dyskowej pozostało itd. Te programy nadzorujące odpowiadają za stan maszyn oraz ich ponowne uruchamianie i rekonfigurowanie przez włączanie, zatrzymywanie i restartowanie różnych programów. Okazuje się, że takie środowisko jest skomplikowane już na jednym serwerze. Pomnóż teraz tę złożoność razy 10 000 serwerów w pomieszczeniu.

1.3. Oprogramowanie w centrum danych

Oprogramowanie dla centrów danych różni się od samodzielnych programów i oprogramowania testowego. Składa się z *wielu warstw* podsystemów, z których liczne działają równolegle i obsługują żądania od wielu różnych usług i licznych instancji pojedynczych usług. Każdy podsystem próbuje zwrócić odpowiedź wystarczająco szybko, aby osiągnąć stawiane mu cele z obszaru latencji. Często się zdarza, że warstwy obsługujące jedno żądanie użytkownika działają na różnych serwerach. Aby zwiększyć wydajność, w wielu warstwach używana jest programowa pamięć podręczna do przechowywania niedawno potrzebnych danych

lub obliczonych wyników. Znalezienie i ponowne użycie danych zapisanych w pamięci podręcznej nazywane jest *trafieniem*; z kolei brak danych w pamięci podręcznej to *chybienie*. Dalej zobaczysz, że dynamika działania pamięci podręcznych może w zaskakujący sposób wpływać na latencję transakcji.

Gdy użytkownik żąda tekstu e-maila, żądanie najpierw trzeba skierować do centrum danych zawierającego główne repozytorium poczty elektronicznej tego użytkownika, a następnie przekazać do serwera równoważącego obciążenie, który to serwer przekieruje je do którejś z mniej obciążonych maszyn spośród setek serwerów odpowiedzialnych za frontend poczty elektronicznej. Warstwa frontendlu zarządza żądaniem i ostatecznie generuje wynik w formacie HTML lub zgodny z API aplikacji. Ta warstwa żąda danego e-maila od warstwy backendu, która kieruje wywołanie do *warstwy bazy danych*, ta wywołuje warstwę pamięci podręcznej bazy danych, a jeśli nastąpi chybienie, kieruje żądanie do warstwy replikacji (aby uzyskać dostęp do rezerwowego repozytorium poczty elektronicznej z innego centrum danych). Ostatecznie zgłaszane jest wywołanie do warstwy serwera dysków, który wczytuje e-mail z jednego z kilku redundantnych dysków, co ilustruje rysunek 1.1. Następnie wyniki są zwracane przez całe drzewo wywołań i często są przy tym modyfikowane.

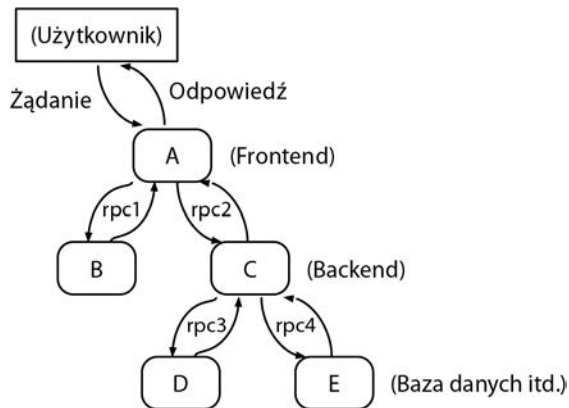


Rysunek 1.1. Warstwy przykładowego oprogramowania

Wszystkie te operacje są łączone za pomocą jakiejś formy przekazywania komunikatów lub zdalnych wywołań procedur (ang. *remote procedure call* — RPC). W tej książce używam skrótu „RPC”. Wywołania RPC z jednej warstwy do drugiej mogą być przetwarzane synchronicznie (jednostka wywołująca oczekuje wtedy na odpowiedź) lub asynchronicznie (jednostka wywołująca kontynuuje pracę i może zgłaszać inne wywołania RPC wykonywane równoległe na wielu różnych serwerach). To właśnie równoległe wykonywanie niewielkich porcji pracy umożliwia oprogramowaniu w centrum danych wykonywanie dla jednego żądania olbrzymich ilości zadań i kończenie ich w ułamku sekundy. Jedna transakcja z udziałem użytkownika może skutkować użyciem od 200 do 2000 różnych serwerów.

Na rysunku 1.2 pokazane jest niewielkie przykładowe drzewo wywołań RPC w stylu z pracy [Sigelman 2010]. Serwer A kieruje synchroniczne wywołanie do serwera B, a po zwróceniu sterowania przez B serwer A kieruje wywołanie do serwera C. Serwer C może równoległe wywołać serwery D i E, a następnie oczekiwać na zwrócenie przez nie odpowiedzi.

Każde żądanie od użytkownika i każde podżądanie RPC mają cel w postaci określonego czasu odpowiedzi. Jeśli dla żądania od użytkownika do warstwy frontendlu poczty elektronicznej celem jest 200 milisekund, dla warstwy backendu cel może wynosić 160 milisekund, dla warstwy



Rysunek 1.2. Drzewo wywołań RPC dla pięciu serwerów (A – E)

bazy danych może być jeszcze krótszy itd., aż do 50 milisekund dla serwera dysków. Gdy podżądanie zbyt długo generuje odpowiedź, każda jednostka wywołująca też jest narażona na przekroczenie czasu. Dla zestawu równoległych wywołań używane jest pojęcie *asymetria czasów wykonania* (ang. *execution skew*) określające zmienność czasu wykonywania zadań.

Jeśli równolegle przetwarzanych jest wiele wywołań RPC, zwykle to najwolniejsze z nich określa łączny czas odpowiedzi. Dlatego jeżeli wykonujesz równoległe 100 wywołań RPC, to łączny czas odpowiedzi jest zależny od czasu na poziomie percentyla 99%. Asymetria czasów wykonania sprawia, że ważne jest zrozumienie i kontrolowanie długich czasów odpowiedzi.

1.4. Latencja z długiego ogona rozkładu

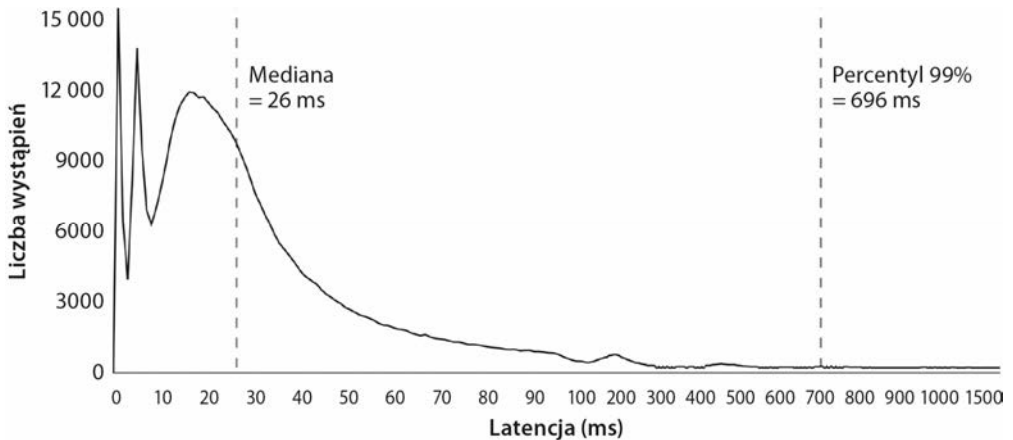
Latencja to czas zegarowy, jaki upływa między dwoma zdarzeniami. W trakcie precyzyjnego omawiania latencji trzeba określić, jakie są to zdarzenia. Na przykład „latencja wywołania RPC” może oznaczać czas między wysłaniem żądania przez program w trybie użytkownika (czyli klienta w terminologii modelu klient-serwer) a otrzymaniem przez niego odpowiedzi na to żądanie. Może też oznaczać czas od otrzymania żądania przez wywołany program w trybie użytkownika (serwer w terminologii modelu klient-serwer) do momentu wysłania odpowiedzi. Latencja w tych dwóch różnych definicjach (dla klienta i dla serwera) dla tego samego wywołania RPC może się różnić o 30 milisekund lub nawet więcej, co może skłaniać do zastanowienia się nad tym, co się stało z dodatkowym czasem i w których urządzeniach komputerowych lub sieciowych on zaginął.

W tej książce skupiam się na latencji wywołań RPC na *serwerze*, chyba że omawiam rozbieżności takie jak 30 milisekund w tym przykładzie.

Żądania RPC kierowane do usługi mają różną latencję, jednak często dla podobnych żądań latencja grupuje się wokół podobnych wartości. Można to zilustrować na histogramie wartości latencji, czyli na wykresie przedstawiającym na osi X przedziały latencji, a na osi Y liczbę wywołań RPC z poszczególnych przedziałów.

W transakcjach w centrum danych histogramy latencji mają jeden lub kilka wierzchołków dla standardowych przypadków i często *długi ogon* znacznie dłuższych czasów odpowiedzi dla nietypowych przypadków [Blake 2015, Hoff 2012, Weaveworks 2017, Dean 2013].

Na histogramie dotyczącym serwera dysków (rysunek 1.3) występują trzy wierzchołki z czasem około 1 milisekundy, 3 milisekund i 20 milisekund dla trzech różnych rodzajów standardowych przypadków, a następnie długi ogon, który ciągnie się do ponad 1500 milisekund. Pożądany czas odpowiedzi wynosi 50 milisekund lub mniej. Ta książka ma pomóc w zrozumieniu i skróceniu długiego ogona. Słabo zaznaczone wierzchołki na poziomie tuż przed 250 milisekundami, 500 milisekundami, 750 milisekundami itd. wskazują na źródło tej zagadki związanej z wydajnością. Jej rozwiązanie znajdziesz w części II.



Rysunek 1.3. Histogram obrazujący dostęp do serwera dysków; po prawej stronie widoczny jest długi ogon

Histogram latencji warto przedstawić za pomocą tylko kilku liczb, zamiast używać 500 wartości, jeśli wyznaczonych jest 500 przedziałów. Jakie wartości należy wybrać?

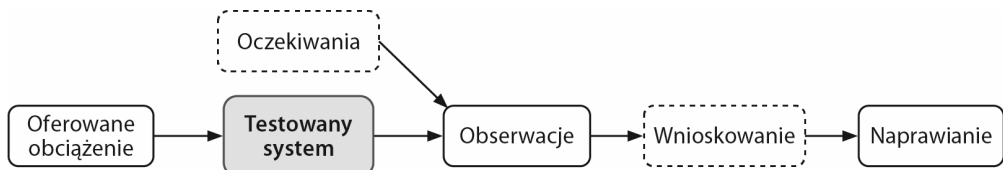
Mediana (lub podobna średnia) latencji jest wyjątkowo słabym opisem asymetrycznego lub wielowierzchołkowego rozkładu, ponieważ rzadko wypada w pobliżu dużej liczby rzeczywistych wartości i nic nie mówi o kształcie ani wielkości długiego ogona, który jest tu interesujący. Na wykresie z rysunku 1.3 mediana wynosi 26 milisekund, co nie niesie żadnych informacji na temat wierzchołków lub ogona. Także maksymalna latencja jest kiepskim wyborem, ponieważ w ciągu dnia może wystąpić jedno skrajnie długie wywołanie RPC (wynikające na przykład z odzyskiwania pamięci lub sprzętowego błędu dysku), podczas gdy wszystkie inne są dziesiątki lub nawet setki razy krótsze.

Zamiast tego lepiej posłużyć się percentylami. Jeśli histogram latencji reprezentuje 50 000 pomiarów, najkrótszych 500 to najszybszy 1%, a najdłuższych 500 to najwolniejszy 1%. Wartość percentyla 99% to liczbowa granica między najszybszymi 99% a najwolniejszym 1%; 99% posortowanych pomiarów jest mniejszych lub równych tej wartości. Istnieje wiele takich wartości. Wszystkie one wypadają między 49 500. a 49 001. posortowanym pomiarem. Można posłużyć się dowolną z tych liczb, jednak zwyczajowo stosuje się 49 500. wartość. Szybkim, ale przydatnym sposobem opisanie rozkładu z długim ogonem jest podanie wartości percentyla 99%, 95%, 99,9% itd. Wartość percentyla 99% w histogramie na rysunku 1.3 wynosi 696 milisekund. Jest ona zbyt wysoka w porównaniu z celem równym 50 milisekund. Oznacza to poważny problem z wydajnością.

Z rozdziału 9. dowiesz się, co spowodowało ten konkretny długi ogon i jak problem został rozwiązany. Poznasz też nową wartość percentyla 99%, która wyniosła około 150 milisekund. Ta prosta zmiana dała zyski równe mojemu 10-letniemu wynagrodzeniu.

1.5. Model myślenia

W procesie analizy latencji z długiego ogona rozkładu i powiązanych problemów z wydajnością stosuję następujące podejście: najpierw należy oszacować, ile czasu jakieś zadanie powinno zająć, następnie trzeba zaobserwować, ile czasu rzeczywiście zajmuje, a na końcu zastanowić się nad różnicami między tymi wartościami. Rysunek 1.4 przedstawia ten model.



Rysunek 1.4. Model analizowania wydajności złożonego oprogramowania

Ten model ujmuje testowany system obejmujący oprogramowanie i sprzęt, oferowane obciążenie tego systemu, oczekiwania co do wydajności systemu przy danym obciążeniu, obserwacje dynamiki działania i wydajności za pomocą narzędzi, wnioskowanie przez człowieka na temat tego, co się dzieje, i ostatecznie wprowadzenie poprawek lub zmian w celu zwiększenia wydajności.

1.6. Szacowanie rzędu wielkości

W trakcie badania wydajności oprogramowania jednym z aspektów pytania „Jak wolne powinno być?” jest oszacowanie, ile czasu powinny zajmować określone zadania. Nawet gdy te szacunki są bardzo ogólne, mogą zapewnić cenne informacje. Programiści zwracający uwagę na wydajność nieustannie szacują w głowach rzędy wielkości, gdy projektują i piszą ważne programy.

Określenie *rzędu wielkości* dotyczy przybliżonej oceny wielkości liczby. Rząd wielkości w systemie *dziesiętnym* wymaga przedstawienia szacunków za pomocą potęg liczby 10 (1, 10, 100 itd.). Gdy określany jest rząd wielkości w systemie *dwójkowym*, należy podać szacunki za pomocą potęg liczby 2 (1, 2, 4, 8 itd.). Możesz się też zetknąć z pośrednimi dziesiętnymi rzędami wielkości — 1, 3, 10, 30 itd. W tej książce używam rzędów wielkości w systemie dziesiętnym, chyba że napiszę, że jest inaczej. Stosuję notację $O(n)$, co oznacza „rzędu n ”, i zawsze podaję *jednostki*. Jest bardzo ważne, czy chodzi o $O(10)$ nanosekund, $O(10)$ milisekund czy $O(10)$ bajtów. Od tego miejsca używam zapisu *ns*, *μs* i *ms* jako skrótów jednostek czasu (nanosekund, mikrosekund i milisekund).

W tabeli 1.1 znajdziesz listę szacunków, które każdy poważny programista powinien znać. Pochodzi ona z wystąpienia Jeffa Deana (jednego z bardzo nielicznych posiadaczy tytułu Google Fellow) z 2009 roku. Od tamtego czasu podane wartości dużo się nie zmieniły. Dodałem tylko kolumnę z rzędami wielkości.

Tabela 1.1. Liczby, które każdy powinien znać [Dean 2009]

Operacja	Czas	O(n)
Dostęp do pamięci podręcznej L1	0,5 ns	O(1) ns
Błędna predykcja gałęzi	5 ns	O(10) ns
Dostęp do pamięci podręcznej L2	7 ns	O(10) ns
Zajęcie lub zwolnienie muteksu	25 ns	O(10) ns
Dostęp do pamięci głównej	100 ns	O(100) ns
Kompresja 1 KB za pomocą biblioteki Zipy	3000 ns	O(1) μ s
Przesyłanie 2 KB przez sieci o szybkości 1 Gb/s	20 000 ns	O(10) μ s
Sekwencyjny odczyt 1 MB danych z pamięci	250 000 ns	O(100) μ s
Przesyłanie pakietu tam i z powrotem w tym samym centrum danych	500 000 ns	O(1) ms
Wyszukiwanie danych na dysku	10 000 000 ns	O(10) ms
Sekwencyjny odczyt 1 MB danych z dysku	20 000 000 ns	O(10) ms
Przesyłanie pakietu na trasie Kalifornia->Holandia->Kalifornia	150 000 000 ns	O(100) ms

Oszacowanie rzędu wielkości oczekiwanego czasu wykonywania różnych części programu sprawia, że po uzyskaniu rzeczywistych pomiarów tych czasów można łatwo dostrzec miejsca, w których wyniki znacznie różnią się od oczekiwań. **Właśnie to umożliwia naukę.** Czasem szacunki są błędne i dowiesz się czegoś nowego na temat działania komputerów lub programów. Czasem szacunki będą poprawne, ale okaże się, że program działa niezgodnie z Twoimi wyobrażeniami, a nieoczekiwane szybsze lub wolniejsze operacje wymagają poprawek. Gdy nabierzesz wprawy i Twoje szacunki staną się trafniejsze, coraz więcej znalezionych rozbieżności będzie oznaczać rzeczywiste błędy z obszaru wydajności.

Znajomość szacunków z tabeli 1.1 pomoże Ci też identyfikować prawdopodobne źródła błędów związanych z wydajnością. Jeśli wykonywanie jakiegoś fragmentu programu zajmuje 100 ms więcej, niż oczekujesz, problem prawdopodobnie nie jest związany z błędną predykcją gałęzi, ponieważ wpływ tego aspektu jest 10 000 000 mniejszy niż 100 ms. Bardziej prawdopodobnym źródłem problemu są operacje dyskowe lub sieciowe albo, o czym przekonasz się w dalszych rozdziałach, długi czas zajmowania blokad lub oczekiwania na długie podżądania RPC.

Dalej zaprojektuj i zbuduj narzędzia do obserwacji i wyświetlania danych. Gdy będziesz z nich korzystać, postaraj się wyrobić sobie nawyk prognozowania (w formie rzędu wielkości), czego spodziewasz się zobaczyć. Po nabraniu wprawy w realizowaniu cyklu predykcje-obszerności-porównania błyskawicznie zaczniesz dostrzegać nietypowe wyniki.

1.7. Dlaczego transakcje działają powoli?

Warto przypomnieć, że szczególnie interesujące są transakcje, które zwykle działają szybko, ale w niektórych sytuacjach transakcje wymagają dużo dodatkowego czasu — na tyle dużo, że użytkownik końcowy odczuwa irytujące opóźnienie. Powolne transakcje to przede wszystkim te,

które naruszają zapisany cel związany z czasem odpowiedzi. Co może powodować takie opóźnienia? Co może być przyczyną zmiennych latencji, a przede wszystkim długiego ogona w ich rozkładzie?

Wskazówką jest to, że dana transakcja (lub określony typ transakcji) jest zwykle szybka. Gdy jest powolna, mamy do czynienia ze standardowym czasem wykonywania zwiększonym o jakieś nieznane opóźnienie. Jeśli uda się zidentyfikować źródło tego opóźnienia, zwykle można wprowadzić proste zmiany w kodzie, które eliminują większość tego opóźnienia, i w ten sposób skrócić długi ogon w rozkładzie latencji.

W oprogramowaniu z wieloma warstwami najczęstszym źródłem opóźnień w jednej warstwie jest oczekiwanie na odpowiedź z niższej warstwy. Najniższa wolna warstwa może działać długo z powodu własnej charakterystyki lub przeciążenia *przesadnym oferowanym obciążeniem*. W trakcie wyznaczania celów dotyczących czasu reakcji koniecznie pamiętaj o tym, aby jednocześnie określić cele (a precyzyjniej — *ograniczenia*) związane z oferowanym obciążeniem.

Zmiana kodu w jednej warstwie nie pomoże, jeśli czeka ona na niższą warstwę. Trzeba znaleźć najniższą warstwę, która jest zbyt wolna, i starać się poprawić jej działanie. Dlatego warto projektować oprogramowanie w taki sposób, aby umożliwić obserwowanie czasu działania każdej warstwy i przekształcać uzyskane pomiary na postać pozwalającą szybko dostrzec, co jest wąskim gardłem. Prosty wykres tego typu jest zestawienie rzeczywistego oferowanego obciążenia z oczekiwanym oferowanym obciążeniem i rzeczywistego czasu odpowiedzi z oczekiwanym w każdej warstwie lub w każdym interfejsie wywołań RPC.

Jeśli oferowane obciążenie w warstwie N jest akceptowalne i warstwa ta nie oczekuje zbyt długo na warstwę niższego poziomu (N + 1), ale i tak czas odpowiedzi warstwy N jest zbyt długi, istnieją wykonywane na jednym serwerze wywołania RPC do warstwy N, których latencja często jest normalna, ale czasem znacznie dłuższa. Warto wtedy dokładnie przyjrzeć się danemu serwerowi. Możliwe, że powolne wywołania RPC wykonują *dotatkową pracę*, która normalnie nie jest realizowana, albo przetwarzają standardowe zadania, ale *działają powoli*, wolniej niż zwykle.

Dodatkowa praca wynika ze struktury gałęzi kodu i przechowywanego stanu. Programy bardzo się różnią między sobą, jeśli chodzi o to, co może być powodem wykonywania dodatkowych zadań. Jednak zwykle realizują dodatkową pracę nawet wtedy, jeśli program działa na serwerze sam, bez innych aplikacji. Tego typu błędy dotyczące wydajności są stosunkowo łatwe do wykrycia. Wystarczy uruchomić kod w trybie offline w środowisku testowym i przekazać do niego sklonowane lub zarejestrowane żądania generowane w związku z rzeczywistym obciążeniem, a przy tym zastosować dodatkową instrumentację w celu znalezienia wzorców rozgałęzień prowadzących do problemu. Wykonywanie kodu na maszynach testowych pozwala zastosować standardowe narzędzia do pomiaru wydajności, które spowalniają przetwarzanie 2-krotnie, a nawet 20-krotnie lub jeszcze bardziej. W książce Brendana Gregga [Gregg 2021] opisanych jest wiele narzędzi do obserwacji odpowiednich do użytku w takim środowisku.

Bardziej ciekawym przypadkiem (i tematem tej książki) są wywołania RPC, które wykonują normalne zadania, ale wolniej niż zazwyczaj. Oznacza to, że coś *oddziałuje* na normalną pracę wywołania RPC na jednym serwerze i spowalnia wykonywanie kodu. Takie transakcje nazywam *hamowanymi*. Ich opóźnienie zwykle nie ujawnia się w testach w trybie offline, a tylko w obliczu rzeczywistego obciążenia generowanego przez użytkowników i często wyłącznie w godzinie podczas dnia, gdy aktywność jest największa. Celem jest wykrycie źródeł oddziaływań i ich wyeliminowanie lub przynajmniej zminimalizowanie ich wpływu. Niestety, w środowisku

produkcyjnym narzędzia do obserwacji, które spowalniają przetwarzanie 2-krotnie czy nawet o 10%, z nadto wydłużają czas odpowiedzi, by można je zastosować. W aktywnych centrach danych (lub w pojazdach, grach wieloosobowych itd.) potrzebne są techniki i narzędzia do obserwacji z kosztami dodatkowymi mniejszymi niż 1%. W branży jest bardzo mało tego typu produktów. Jeden z nich poznasz w części III tej książki.

Warto przypomnieć, że w centrum danych na każdym serwerze działa wiele programów, a każdy z nich może mieć wiele wątków. Oddziaływania na jednym serwerze muszą być spowodowane przez jakieś czynniki z tej maszyny (może to być przychodzący i wychodzący ruch sieciowy). Oddziaływania w takim środowisku prawie zawsze są wynikiem rywalizacji o współdzielony zasób.

1.8. Pięć podstawowych zasobów

Są tylko cztery zasoby sprzętowe współdzielone przez niepowiązane programy działające na jednym serwerze:

- procesor,
- pamięć,
- dysk twardy lub nośnik SSD,
- sieć.

Jeśli w programie działa kilka współpracujących ze sobą wątków, występuje też piąty podstawowy zasób:

- sekcja krytyczna w oprogramowaniu.

Sekcja krytyczna to fragment kodu używający współdzielonych danych w sposób, który może skutkować nieprawidłowym działaniem programu, jeśli ten fragment będzie wykonywany przez więcej niż jeden wątek naraz. Tego rodzaju sekcje kodu są chronione za pomocą blokad programowych, aby jednocześnie mogły być wykonywane tylko przez jeden wątek. Inne wątki muszą wtedy czekać na wejście do sekcji krytycznej.

Aby wykryć oddziaływania, trzeba najpierw zrozumieć normalne działanie kodu. Punktem wyjścia jest nauczenie się technik precyzyjnego pomiaru pięciu podstawowych zasobów. W pozostałych rozdziałach części I omawiam cztery zasoby sprzętowe, a blokady programowe opisuję dopiero w rozdziale 27., kiedy odpowiednie narzędzie do obserwacji będzie już gotowe. Jeśli powolne wywołanie RPC próbuje użyć jednego z pięciu wymienionych zasobów, a inny program lub wątek także korzysta z danego zasobu, wywołanie RPC będzie musiało na niego poczekać. Jest to główny mechanizm będący źródłem oddziaływań lub hamowania.

1.9. Podsumowanie

Ta książka ma pomóc w zrozumieniu dynamiki działania transakcji w oprogramowaniu centrów danych, baz danych, komputerów stacjonarnych, gier i dedykowanych kontrolerów, a przede wszystkim tych transakcji, których działanie zajmuje znacznie więcej czasu niż zazwyczaj.

Dobry programista potrafi oszacować z precyzją na poziomie rzędu wielkości, jak dużo czasu powinno zająć wykonywanie danego fragmentu kodu. Dzięki temu umie też wykryć i poprawić kod, który zawsze działa zbyt wolno. W tej książce zakładam, że kod, który zawsze jest wykonywany zbyt długo, został już poprawiony. Interesuje nas znacznie trudniejszy do zrozumienia kod, który tylko czasem działa zbyt powoli.

Na serwerze w centrum danych, który przetwarza tysiące transakcji na sekundę, niektóre z nich będą niekiedy wykonywane powoli, a następnie ponownie szybko. Na histogramie czasów transakcji widać długi ogon wolnych przebiegów, które mają nieproporcjonalny wpływ na łączny czas odpowiedzi dla użytkowników i nieproporcjonalnie zmniejszają ilość pracy, jaką dany serwer potrafi wykonać. Te powolne transakcje są wynikiem jakiegoś rodzaju oddziaływań, jednak w wielowarstwowym oprogramowaniu w centrum danych często trudno jest ustalić, która warstwa jest powolna w określonej transakcji. Dlatego trudno też jest stwierdzić, gdzie należy szukać oddziaływań.

Szacunki rzędu wielkości, takie jak pokazane w tabeli 1.1, mogą być wskazówką przy identyfikowaniu możliwych źródeł lub mechanizmów związanych z problemami z wydajnością, jednak przeważnie nie da się precyzyjnie ustalić powolnego fragmentu kodu. Aby to zrobić, trzeba zaprojektować odpowiednie narzędzia do obserwacji warstwowego oprogramowania i serwerów wykonujących wiele niepowiązanych programów, które mogą wpływać na swoje działanie.

Na ogólnym poziomie transakcja na pojedynczym serwerze albo działa normalnie, albo działa powoli, albo oczekuje na coś z tego serwera. Dwa ostatnie z tych zjawisk są skutkiem interakcji. Dalej omawiam mechanizmy związane z tymi zjawiskami i pokazuję, jak zaobserwować je *in situ*.

I to tyle. Aby rozwiązać problemy z okazjonalnie powolnymi transakcjami, wystarczy: (i) zidentyfikować, która warstwa kodu działa powoli, a następnie (ii) ustalić, co wpływa na jej działanie, a potem (iii) wprowadzić poprawki. Reszta książki ma nauczyć Cię, jak wykonywać te trzy proste kroki. Niestety dwa pierwsze z nich nie są łatwe.

- Skupiam się na zrozumieniu okazjonalnie powolnych transakcji RPC.
- Gdy równolegle wykonywanych jest 100 wywołań RPC, czas z poziomu percentyla 99% jest wyznacznikiem *ogólnego* czasu odpowiedzi.
- Oprogramowanie w centrach danych jest podatne na asymetryczny rozkład czasu wykonywania, w którym występuje długi ogon bardzo wolnych odpowiedzi.
- Spowolnione lub hamowane transakcje oznaczają, że coś oddziałuje na pracę wywołań RPC.
- Oddziaływania wynikają ze współdzielenia pięciu podstawowych zasobów.
- Oddziaływania są trudne do zaobserwowania *in situ*. Dalej zbuduję brakujące narzędzia do obserwacji.
- Koniecznie oszacuj rząd wielkości oczekiwanych czasów wykonywania, aby ułatwić sobie dostrzeżenie nieoczekiwanych wyników.

PROGRAM PARTNERSKI

— GRUPY HELION —



1. ZAREJESTRUJ SIĘ
2. PREZENTUJ KSIĄŻKI
3. ZBIERAJ PROWIZJĘ

Zmień swoją stronę WWW w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

<http://program-partnerski.helion.pl>

GRUPA
Helion 

Kod eksperta: tu każdy bit ma znaczenie!

Mimo że dysponujemy coraz szybszym i wydajniejszym sprzętem, oprogramowaniu wciąż stawia się wysrubowane wymagania: musi działać błyskawicznie i niezawodnie. Nierzadko od jakości pracy systemu zależy bardzo dużo, nawet bezpieczeństwo i życie człowieka. Eliminacja źródłowych przyczyn problemów wydajnościowych jest jednak niezwykle trudnym zadaniem. Wymaga wiedzy o działaniu sprzętu, interakcjach między komponentami systemu, a także wysokiej klasy umiejętności programistycznych.

Ta książka jest przeznaczona dla programistów i osób zaawansowanych w nauce programowania. Dzięki niej nauczysz się uwzględniać w projektach narzędzia do obserwacji pracy kodu i analizować uzyskane dzięki nim dane o wydajności. Dowiesz się też, jak na podstawie takiego wnioskowania uzyskiwać znaczną poprawę w szybkości przetwarzania powolnych operacji. Prezentowane treści zilustrowano przykładami i ćwiczeniami w języku C lub C++ w systemie Linux. Materiał zawarty w książce to trzy główne zagadnienia: pomiary, obserwacje, a także wnioskowanie i nanoszenie poprawek w kodzie. Wartościową częścią książki jest omówienie procesu budowy niskokosztowego narzędzia do obserwacji KTrace i jego zastosowania we wdrażaniu wyrafinowanych rozwiązań programistycznych. W ten sposób można u źródła usuwać przyczyny problemów z wydajnością kodu.

Najciekawsze zagadnienia:

- problemy sprzętowe: procesory, pamięci, dyski twarde, nośniki SSD i sieci
- korygowanie kodu wolno działającego programu
- przydatne mechanizmy monitorowania pracy kodu
- analiza danych dotyczących wydajności
- identyfikacja problemów, takich jak wykonywanie zbyt wielu instrukcji, powolne wykonywanie instrukcji, oczekiwanie na zasoby i blokady programowe

DR RICHARD L. SITES zajmuje się programowaniem od 1959 roku. Zdobył olbrzymią wiedzę o interakcjach sprzętu i oprogramowania. Rozwijał mikrokod architektury VAX, był jednym z architektów mikroprocesora DEC Alpha i wymyślił liczniki wydajności powszechnie stosowane w procesorach. Zajmował się niskokosztowym śledzeniem mikrokodu i oprogramowania w firmach: DEC, Adobe, Google i Tesla. Posiada 66 patentów i jest członkiem National Academy of Engineering.

Helion 

 **helion.pl**

 **HELION SA**
ul. Kościuszki 1c
44-100 Gliwice
tel.: 32 230 98 63
helion@helion.pl

KOD KORZYŚCI
Sięgnij po więcej! ▶



ISBN 978-83-283-9515-2



9 788328 395152

Cena: 119,00 zł

 **Pearson**
Addison-Wesley