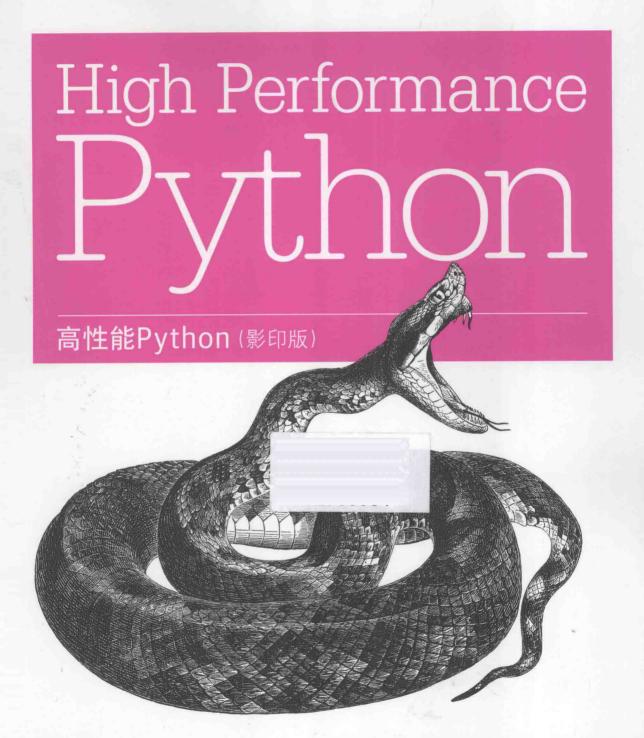
O'REILLY



高性能Python (影印版) **High Performance Python**

Micha Gorelick, Ian Ozsvald 著

Beijing · Cambridge · Farnham · Köln · Sebastopol · Tokyo O'REILLY®



O'Reilly Media, Inc.授权东南大学出版社出版

南京 东南大学出版社

图书在版编目(CIP)数据

高性能 Python:英文/(美)戈雷利克(Gorelick, M.), (英)欧日沃尔德(Ozsvald, I.)著. 一影印本. 一南京:东南大学出版社,2015.2

书名原文: High Performance Python ISBN 978-7-5641-5385-4

I.①高··· II.①戈··· ②欧··· III.①软件工具—程序设计—英文 IV. ①TP311.56

中国版本图书馆 CIP 数据核字(2014)第 294385 号 江苏省版权局著作权合同登记 图字:10-2013-366号

© 2014 by O'Reilly Media, Inc.

Reprint of the English Edition, iointly published by O'Reilly Media, Inc. and Southeast University Press, 2015. Authorized reprint of the original English edition, 2013 O'Reilly Media, Inc., the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 O'Reilly Media, Inc.出版 2014。

英文影印版由东南大学出版社出版 2015。此影印版的出版和销售得到出版权和销售权的所有者—— O'Reilly Media, Inc.的许可。

版权所有,未得书面许可,本书的任何部分和全部不得以任何形式重制。

高性能 Pvthon(影印版)

出版发行: 东南大学出版社

地 址:南京四牌楼 2号 邮编:210096

出版人: 江建中

网 址: http://www.seupress.com

电子邮件: press@seupress.com

印 刷:常州市武进第三印刷有限公司

开 本: 787 毫米×980 毫米 16 开本

印 张: 23.25

字 数: 455 千字

版 次:2015年2月第1版

印 次: 2015年2月第1次印刷

书 号: ISBN 978-7-5641-5385-4

定 价: 78.00元

Preface

Python is easy to learn. You're probably here because now that your code runs correctly, you need it to run faster. You like the fact that your code is easy to modify and you can iterate with ideas quickly. The trade-off between *easy to develop* and *runs as quickly as I need* is a well-understood and often-bemoaned phenomenon. There are solutions.

Some people have serial processes that have to run faster. Others have problems that could take advantage of multicore architectures, clusters, or graphics processing units. Some need scalable systems that can process more or less as expediency and funds allow, without losing reliability. Others will realize that their coding techniques, often borrowed from other languages, perhaps aren't as natural as examples they see from others.

In this book we will cover all of these topics, giving practical guidance for understanding bottlenecks and producing faster and more scalable solutions. We also include some war stories from those who went ahead of you, who took the knocks so you don't have to.

Python is well suited for rapid development, production deployments, and scalable systems. The ecosystem is full of people who are working to make it scale on your behalf, leaving you more time to focus on the more challenging tasks around you.

Who This Book Is For

You've used Python for long enough to have an idea about why certain things are slow and to have seen technologies like Cython, numpy, and PyPy being discussed as possible solutions. You might also have programmed with other languages and so know that there's more than one way to solve a performance problem.

While this book is primarily aimed at people with CPU-bound problems, we also look at data transfer and memory-bound solutions. Typically these problems are faced by scientists, engineers, quants, and academics.

We also look at problems that a web developer might face, including the movement of data and the use of just-in-time (JIT) compilers like PyPy for easy-win performance gains.

It might help if you have a background in C (or C++, or maybe Java), but it isn't a prerequisite. Python's most common interpreter (CPython-the standard you normally get if you type python at the command line) is written in C, and so the hooks and libraries all expose the gory inner C machinery. There are lots of other techniques that we cover that don't assume any knowledge of C.

You might also have a lower-level knowledge of the CPU, memory architecture, and data buses, but again, that's not strictly necessary.

Who This Book Is Not For

This book is meant for intermediate to advanced Python programmers. Motivated novice Python programmers may be able to follow along as well, but we recommend having a solid Python foundation.

We don't cover storage-system optimization. If you have a SQL or NoSQL bottleneck, then this book probably won't help you.

What You'll Learn

Your authors have been working with large volumes of data, a requirement for *I want* the answers faster! and a need for scalable architectures, for many years in both industry and academia. We'll try to impart our hard-won experience to save you from making the mistakes that we've made.

At the start of each chapter, we'll list questions that the following text should answer (if it doesn't, tell us and we'll fix it in the next revision!).

We cover the following topics:

- Background on the machinery of a computer so you know what's happening behind the scenes
- Lists and tuples—the subtle semantic and speed differences in these fundamental data structures
- Dictionaries and sets—memory allocation strategies and access algorithms in these important data structures
- Iterators—how to write in a more Pythonic way and open the door to infinite data streams using iteration
- · Pure Python approaches—how to use Python and its modules effectively

- Matrices with numpy—how to use the beloved numpy library like a beast
- Compilation and just-in-time computing—processing faster by compiling down to machine code, making sure you're guided by the results of profiling
- · Concurrency—ways to move data efficiently
- multiprocessing—the various ways to use the built-in multiprocessing library for parallel computing, efficiently share numpy matrices, and some costs and benefits of interprocess communication (IPC)
- Cluster computing—convert your multiprocessing code to run on a local or remote cluster for both research and production systems
- Using less RAM—approaches to solving large problems without buying a humungous computer
- Lessons from the field—lessons encoded in war stories from those who took the blows so you don't have to

Python 2.7

Python 2.7 is the dominant version of Python for scientific and engineering computing. 64-bit is dominant in this field, along with *nix environments (often Linux or Mac). 64bit lets you address larger amounts of RAM. *nix lets you build applications that can be deployed and configured in well-understood ways with well-understood behaviors.

If you're a Windows user, then you'll have to buckle up. Most of what we show will work just fine, but some things are OS-specific, and you'll have to research a Windows solution. The biggest difficulty a Windows user might face is the installation of modules: research in sites like StackOverflow should give you the solutions you need. If you're on Windows, then having a virtual machine (e.g., using VirtualBox) with a running Linux installation might help you to experiment more freely.

Windows users should definitely look at a packaged solution like those available through Anaconda, Canopy, Python(x,y), or Sage. These same distributions will make the lives of Linux and Mac users far simpler too.

Moving to Python 3

Python 3 is the future of Python, and everyone is moving toward it. Python 2.7 will nonetheless be around for many years to come (some installations still use Python 2.4 from 2004); its retirement date has been set at 2020.

The shift to Python 3.3+ has caused enough headaches for library developers that people have been slow to port their code (with good reason), and therefore people have been slow to adopt Python 3. This is mainly due to the complexities of switching from a mix of string and Unicode datatypes in complicated applications to the Unicode and byte implementation in Python 3.

Typically, when you want reproducible results based on a set of trusted libraries, you don't want to be at the bleeding edge. High performance Python developers are likely to be using and trusting Python 2.7 for years to come.

Most of the code in this book will run with little alteration for Python 3.3+ (the most significant change will be with print turning from a statement into a function). In a few places we specifically look at improvements that Python 3.3+ provides. One item that might catch you out is the fact that / means integer division in Python 2.7, but it becomes float division in Python 3. Of course—being a good developer, your wellconstructed unit test suite will already be testing your important code paths, so you'll be alerted by your unit tests if this needs to be addressed in your code.

scipy and numpy have been Python 3-compatible since late 2010. matplotlib was compatible from 2012, scikit-learn was compatible in 2013, and NLTK is expected to be compatible in 2014. Django has been compatible since 2013. The transition notes for each are available in their repositories and newsgroups; it is worth reviewing the processes they used if you're going to migrate older code to Python 3.

We encourage you to experiment with Python 3.3+ for new projects, but to be cautious with libraries that have only recently been ported and have few users—you'll have a harder time tracking down bugs. It would be wise to make your code Python 3.3+compatible (learn about the __future__ imports), so a future upgrade will be easier.

Two good guides are "Porting Python 2 Code to Python 3" (http://bit.ly/pyporting) and "Porting to Python 3: An in-depth guide." (http://python3porting.com/) With a distribution like Anaconda or Canopy, you can run both Python 2 and Python 3 simultaneously—this will simplify your porting.

License

This book is licensed under Creative Commons Attribution-NonCommercial-NoDerivs 3.0 (http://bit.ly/CC_A-NC-ND3).

You're welcome to use this book for noncommercial purposes, including for noncommercial teaching. The license only allows for complete reproductions; for partial reproductions, please contact O'Reilly (see "How to Contact Us" on page xv). Please attribute the book as noted in the following section.

We negotiated that the book should have a Creative Commons license so the contents could spread further around the world. We'd be quite happy to receive a beer if this decision has helped you. We suspect that the O'Reilly staff would feel similarly about the beer.

How to Make an Attribution

The Creative Commons license requires that you attribute your use of a part of this book. Attribution just means that you should write something that someone else can follow to find this book. The following would be sensible: "High Performance Python by Micha Gorelick and Ian Ozsvald (O'Reilly). Copyright 2014 Micha Gorelick and Ian Ozsvald, 978-1-449-36159-4."

Errata and Feedback

We encourage you to review this book on public sites like Amazon—please help others understand if they'd benefit from this book! You can also email us at *feedback@highper formancepython.com*.

We're particularly keen to hear about errors in the book, successful use cases where the book has helped you, and high performance techniques that we should cover in the next edition. You can access the page for this book at http://bit.ly/High_Performance_Python.

Complaints are welcomed through the instant-complaint-transmission-service > /dev/null.

Conventions Used in This Book

The following typographical conventions are used in this book:

Italic

Indicates new terms, URLs, email addresses, filenames, and file extensions.

Constant width

Used for program listings, as well as within paragraphs to refer to commands, modules, and program elements such as variable or function names, databases, datatypes, environment variables, statements, and keywords.

Constant width bold

Shows commands or other text that should be typed literally by the user.

Constant width italic

Shows text that should be replaced with user-supplied values or by values determined by context.



This element signifies a question or exercise.



This element signifies a general note.



This element indicates a warning or caution.

Using Code Examples

Supplemental material (code examples, exercises, etc.) is available for download at https://github.com/mynameisfiber/high_performance_python.

This book is here to help you get your job done. In general, if example code is offered with this book, you may use it in your programs and documentation. You do not need to contact us for permission unless you're reproducing a significant portion of the code. For example, writing a program that uses several chunks of code from this book does not require permission. Selling or distributing a CD-ROM of examples from O'Reilly books does require permission. Answering a question by citing this book and quoting example code does not require permission. Incorporating a significant amount of example code from this book into your product's documentation does require permission.

If you feel your use of code examples falls outside fair use or the permission given above, feel free to contact us at permissions@oreilly.com.

Safari® Books Online



Safari Books Online is an on-demand digital library that delivers expert content in both book and video form from the world's leading authors in technology and business.

Technology professionals, software developers, web designers, and business and creative professionals use Safari Books Online as their primary resource for research, problem solving, learning, and certification training.

Safari Books Online offers a range of plans and pricing for enterprise, government, education, and individuals.

Members have access to thousands of books, training videos, and prepublication manuscripts in one fully searchable database from publishers like O'Reilly Media, Prentice Hall Professional, Addison-Wesley Professional, Microsoft Press, Sams, Que, Peachpit Press, Focal Press, Cisco Press, John Wiley & Sons, Syngress, Morgan Kaufmann, IBM

Redbooks, Packt, Adobe Press, FT Press, Apress, Manning, New Riders, McGraw-Hill, Jones & Bartlett, Course Technology, and hundreds more. For more information about Safari Books Online, please visit us online.

How to Contact Us

Please address comments and questions concerning this book to the publisher:

O'Reilly Media, Inc. 1005 Gravenstein Highway North Sebastopol, CA 95472 800-998-9938 (in the United States or Canada) 707-829-0515 (international or local) 707-829-0104 (fax)

To comment or ask technical questions about this book, send email to bookques tions@oreilly.com.

For more information about our books, courses, conferences, and news, see our website at http://www.oreilly.com.

Find us on Facebook: http://facebook.com/oreilly

Follow us on Twitter: http://twitter.com/oreillymedia

Watch us on YouTube: http://www.youtube.com/oreillymedia

Acknowledgments

Thanks to Jake Vanderplas, Brian Granger, Dan Foreman-Mackey, Kyran Dale, John Montgomery, Jamie Matthews, Calvin Giles, William Winter, Christian Schou Oxvig, Balthazar Rouberol, Matt "snakes" Reiferson, Patrick Cooper, and Michael Skirpan for invaluable feedback and contributions. Ian thanks his wife Emily for letting him disappear for 10 months to write this (thankfully she's terribly understanding). Micha thanks Elaine and the rest of his friends and family for being so patient while he learned to write. O'Reilly are also rather lovely to work with.

Our contributors for the "Lessons from the Field" chapter very kindly shared their time and hard-won lessons. We give thanks to Ben Jackson, Radim Řehůřek, Sebastjan Trebca, Alex Kelly, Marko Tasic, and Andrew Godwin for their time and effort.

Table of Contents

Pre	face	ix
1.	Understanding Performant Python	1
	The Fundamental Computer System	1
	Computing Units	2
	Memory Units	5
	Communications Layers	7
	Putting the Fundamental Elements Together	9
	Idealized Computing Versus the Python Virtual Machine	10
	So Why Use Python?	13
2.	Profiling to Find Bottlenecks	17
	Profiling Efficiently	18
	Introducing the Julia Set	19
	Calculating the Full Julia Set	23
	Simple Approaches to Timing—print and a Decorator	26
	Simple Timing Using the Unix time Command	29
	Using the cProfile Module	31
	Using runsnakerun to Visualize cProfile Output	36
	Using line_profiler for Line-by-Line Measurements	37
	Using memory_profiler to Diagnose Memory Usage	42
	Inspecting Objects on the Heap with heapy	48
	Using dowser for Live Graphing of Instantiated Variables	50
	Using the dis Module to Examine CPython Bytecode	52
	Different Approaches, Different Complexity	54
	Unit Testing During Optimization to Maintain Correctness	56
	No-op @profile Decorator	57
	Strategies to Profile Your Code Successfully	59
	Wrap-Up	60

3.	Lists and Tuples	61
	A More Efficient Search	64
	Lists Versus Tuples	66
	Lists as Dynamic Arrays	67
	Tuples As Static Arrays	70
	Wrap-Up	72
4.	Dictionaries and Sets	. 73
	How Do Dictionaries and Sets Work?	77
	Inserting and Retrieving	77
	Deletion	80
	Resizing	81
	Hash Functions and Entropy	81
	Dictionaries and Namespaces	85
	Wrap-Up	88
5.	Iterators and Generators	89
	Iterators for Infinite Series	92
	Lazy Generator Evaluation	94
	Wrap-Up	98
6.	Matrix and Vector Computation	99
	Introduction to the Problem	100
	Aren't Python Lists Good Enough?	105
	Problems with Allocating Too Much	106
	Memory Fragmentation	109
	Understanding perf	111
	Making Decisions with perf's Output	113
	Enter numpy	114
	Applying numpy to the Diffusion Problem	117
	Memory Allocations and In-Place Operations Selective Optimizations: Finding What Needs to Be Fixed	120 124
	numexpr: Making In-Place Operations Faster and Easier	124
	A Cautionary Tale: Verify "Optimizations" (scipy)	129
	Wrap-Up	130
7	Compiling to C	135
٠.	Compiling to C	136
	JIT Versus AOT Compilers	138
	Why Does Type Information Help the Code Run Faster?	138
	Using a C Compiler	139
	Reviewing the Julia Set Example	140
		140
		1 10

	Compiling a Pure-Python Version Using Cython	141
	Cython Annotations to Analyze a Block of Code	143
	Adding Some Type Annotations	145
	Shed Skin	150
	Building an Extension Module	151
	The Cost of the Memory Copies	153
	Cython and numpy	154
	Parallelizing the Solution with OpenMP on One Machine	155
	Numba	157
	Pythran	159
	PyPy	160
	Garbage Collection Differences	161
	Running PyPy and Installing Modules	162
	When to Use Each Technology	164
	Other Upcoming Projects	165
	A Note on Graphics Processing Units (GPUs)	166
	A Wish for a Future Compiler Project	166
	Foreign Function Interfaces	167
	ctypes	167
	cffi	170
	f2py	173
	CPython Module	175
	Wrap-Up	179
8.	Concurrency	181
	Introduction to Asynchronous Programming	182
	Serial Crawler	185
	gevent	187
	tornado	192
	AsyncIO	196
	Database Example	198
	Wrap-Up	201
€.	The multiprocessing Module	203
	An Overview of the Multiprocessing Module	206
	Estimating Pi Using the Monte Carlo Method	208
	Estimating Pi Using Processes and Threads	210
	Using Python Objects	210
	Using Python Objects Random Numbers in Parallel Systems	210217
	Random Numbers in Parallel Systems Using numpy Finding Prime Numbers	217
	Random Numbers in Parallel Systems Using numpy	217 218

	Serial Solution	236
	Naive Pool Solution	236
	A Less Naive Pool Solution	238
	Using Manager. Value as a Flag	239
	Using Redis as a Flag	241
	Using RawValue as a Flag	243
	Using mmap as a Flag	244
	Using mmap as a Flag Redux	245
	Sharing numpy Data with multiprocessing	248
	Synchronizing File and Variable Access	254
	File Locking	255
	Locking a Value	258
	Wrap-Up	261
10.	Clusters and Job Queues	263
	Benefits of Clustering	264
	Drawbacks of Clustering	265
	\$462 Million Wall Street Loss Through Poor Cluster Upgrade Strategy	266
	Skype's 24-Hour Global Outage	267
	Common Cluster Designs	268
	How to Start a Clustered Solution	268
	Ways to Avoid Pain When Using Clusters	269
	Three Clustering Solutions	270
	Using the Parallel Python Module for Simple Local Clusters	271
	Using IPython Parallel to Support Research	273
	NSQ for Robust Production Clustering	277
	Queues	277
	Pub/sub	278
	Distributed Prime Calculation	280
	Other Clustering Tools to Look At	284
	Wrap-Up	285
11.	Using Less RAM	287
	Objects for Primitives Are Expensive	288
	The Array Module Stores Many Primitive Objects Cheaply	289
	Understanding the RAM Used in a Collection	292
	Bytes Versus Unicode	294
	Efficiently Storing Lots of Text in RAM	295
	Trying These Approaches on 8 Million Tokens	296
	Tips for Using Less RAM	304
	Probabilistic Data Structures	305
	Very Approximate Counting with a 1-byte Morris Counter	306
	K-Minimum Values	308

	Bloom Filters	312
	LogLog Counter	317
	Real-World Example	321
12.	Lessons from the Field	325
	Adaptive Lab's Social Media Analytics (SoMA)	325
	Python at Adaptive Lab	326
	SoMA's Design	326
	Our Development Methodology	327
	Maintaining SoMA	327
	Advice for Fellow Engineers	328
	Making Deep Learning Fly with RadimRehurek.com	328
	The Sweet Spot	328
	Lessons in Optimizing	330
	Wrap-Up	332
	Large-Scale Productionized Machine Learning at Lyst.com	333
	Python's Place at Lyst	333
	Cluster Design	333
	Code Evolution in a Fast-Moving Start-Up	333
	Building the Recommendation Engine	334
	Reporting and Monitoring	334
	Some Advice	335
	Large-Scale Social Media Analysis at Smesh	335
	Python's Role at Smesh	335
	The Platform	336
	High Performance Real-Time String Matching	336
	Reporting, Monitoring, Debugging, and Deployment	338
	PyPy for Successful Web and Data Processing Systems	339
	Prerequisites	339
	The Database	340
	The Web Application	340
	OCR and Translation	341
	Task Distribution and Workers	341
	Conclusion	341
	Task Queues at Lanyrd.com	342
	Python's Role at Lanyrd	342
	Making the Task Queue Performant	343
	Reporting, Monitoring, Debugging, and Deployment	343
	Advice to a Fellow Developer	343
Inc	day.	2/15

Understanding Performant Python

Questions You'll Be Able to Answer After This Chapter

- What are the elements of a computer's architecture?
- What are some common alternate computer architectures?
- How does Python abstract the underlying computer architecture?
- What are some of the hurdles to making performant Python code?
- · What are the different types of performance problems?

Programming computers can be thought of as moving bits of data and transforming them in special ways in order to achieve a particular result. However, these actions have a time cost. Consequently, high performance programming can be thought of as the act of minimizing these operations by either reducing the overhead (i.e., writing more efficient code) or by changing the way that we do these operations in order to make each one more meaningful (i.e., finding a more suitable algorithm).

Let's focus on reducing the overhead in code in order to gain more insight into the actual hardware on which we are moving these bits. This may seem like a futile exercise, since Python works quite hard to abstract away direct interactions with the hardware. However, by understanding both the best way that bits can be moved in the real hardware and the ways that Python's abstractions force your bits to move, you can make progress toward writing high performance programs in Python.

The Fundamental Computer System

The underlying components that make up a computer can be simplified into three basic parts: the computing units, the memory units, and the connections between them. In

addition, each of these units has different properties that we can use to understand them. The computational unit has the property of how many computations it can do per second, the memory unit has the properties of how much data it can hold and how fast we can read from and write to it, and finally the connections have the property of how fast they can move data from one place to another.

Using these building blocks, we can talk about a standard workstation at multiple levels of sophistication. For example, the standard workstation can be thought of as having a central processing unit (CPU) as the computational unit, connected to both the random access memory (RAM) and the hard drive as two separate memory units (each having different capacities and read/write speeds), and finally a bus that provides the connections between all of these parts. However, we can also go into more detail and see that the CPU itself has several memory units in it: the L1, L2, and sometimes even the L3 and L4 cache, which have small capacities but very fast speeds (from several kilobytes to a dozen megabytes). These extra memory units are connected to the CPU with a special bus called the *backside bus*. Furthermore, new computer architectures generally come with new configurations (for example, Intel's Nehalem CPUs replaced the front-side bus with the Intel QuickPath Interconnect and restructured many connections). Finally, in both of these approximations of a workstation we have neglected the network connection, which is effectively a very slow connection to potentially many other computing and memory units!

To help untangle these various intricacies, let's go over a brief description of these fundamental blocks.

Computing Units

The computing unit of a computer is the centerpiece of its usefulness—it provides the ability to transform any bits it receives into other bits or to change the state of the current process. CPUs are the most commonly used computing unit; however, graphics processing units (GPUs), which were originally typically used to speed up computer graphics but are becoming more applicable for numerical applications, are gaining popularity due to their intrinsically parallel nature, which allows many calculations to happen simultaneously. Regardless of its type, a computing unit takes in a series of bits (for example, bits representing numbers) and outputs another set of bits (for example, representing the sum of those numbers). In addition to the basic arithmetic operations on integers and real numbers and bitwise operations on binary numbers, some computing units also provide very specialized operations, such as the "fused multiply add" operation, which takes in three numbers, A,B,C, and returns the value A * B + C.

The main properties of interest in a computing unit are the number of operations it can do in one cycle and how many cycles it can do in one second. The first value is measured