

JULIEN DANJOU

THE HACKER'S GUIDE TO SCALING PYTHON



© 2017 Julien Danjou. All rights reserved.

ISBN 978-1-387-37932-3

Contents

1	Scaling?	14
1.1	Across CPUs	16
1.2	Distributed Systems	19
1.3	Service-Oriented Architecture	20
2	CPU Scaling	23
2.1	Using Threads	23
2.2	Using Processes	27
2.3	Using Futures	30
2.4	Advanced Futures Usage	35
2.5	Daemon Processes	39
2.6	Mehdi Abaakouk on CPU Scaling	46
3	Event Loops	53
3.1	Basic Pattern	54
3.2	Using Asyncio	56
3.2.1	Network Server	63
3.3	Naoki Inada on asyncio	70

4	Functional Programming	75
4.1	The Functional Toolkit	77
5	Queue-Based Distribution	87
5.1	RQ	89
5.2	Celery	94
5.2.1	Handling Failures	96
5.2.2	Chaining Tasks	98
5.2.3	Multiple Queues	99
5.2.4	Monitoring	101
5.3	Joshua Harlow on Task Distribution	103
6	Designing for Failure	111
6.1	Naive Retrying	112
6.2	Retrying with <i>Tenacity</i>	114
7	Lock Management	121
7.1	Thread Locks	122
7.2	Processes Locks	126
7.2.1	Multiprocessing Locks	126
7.2.2	Inter-Processes Locks	129
7.3	Using <i>etcd</i> for Distributed Locking	131
7.4	Using <i>Tooz</i> Locking Abstraction	135

8	Group membership	140
8.1	Creating, Joining and Leaving Groups	140
8.2	Using Capabilities	143
8.3	Using Watchers Callbacks	146
8.4	Consistent Hash Rings	149
8.5	Partitioner	155
8.6	Alexys Jacob-Monier on Cluster Management	161
9	Building REST API	171
9.1	The WSGI Protocol	173
9.2	Streaming Data	177
9.3	Using ETag	183
9.4	Asynchronous HTTP API	190
9.5	Fast HTTP Client	194
9.6	Testing REST API	203
9.7	Chris Dent on HTTP	207
10	Deploying on <i>PaaS</i>	216
10.1	Heroku	217
10.2	Amazon Beanstalk	221
10.3	Google App Engine	224
10.4	OpenShift	227
10.5	Beyond PaaS	232

11 Testing Distributed Systems	234
11.1 Setting Up Environments with <i>tox</i>	235
11.2 Manage External Services with <i>pipfaf</i>	239
11.3 Using Fixtures with <i>pipfaf</i>	244
12 Caching	250
12.1 Local Caching	251
12.2 Memoization	254
12.3 Distributed Caching	257
12.4 Jason Myers on Databases	262
13 Performance	267
13.1 Memory and CPU Profiling	268
13.2 Profiling Strategy and a Case	272
13.3 Zero-Copy	280
13.4 Disassembling Code	287
13.5 Victor Stinner on Performance	292

List of Figures

2.1	Using threads with CPython	27
2.2	Using processes with CPython	29
3.1	<code>asyncio.Protocol</code> state machine	66
5.1	Queue architecture	88
5.2	rq-dashboard example	93
5.3	Celery Flower dashboard	103
7.1	Simple <i>etcd</i> lock algorithm	132
8.1	<i>Tooz</i> membership	141
8.2	Consistent hash ring	150
10.1	Heroku Web dashboard	221
10.2	AWS Elastic Beanstalk dashboard	224
10.3	Google App Engine dashboard	226
10.4	OpenShift Web dashboard	232
13.1	<i>KCacheGrind</i> example	270
13.2	<i>RunSnake</i> example	271

LIST OF FIGURES

13.3 Carbonara profiling information before optimization	273
13.4 Carbonara call graph before optimization	275
13.5 Carbonara profiling information after optimization	277
13.6 Carbonara call graph after optimization	279

List of Examples

1.1	Thread-unsafe code without the GIL	17
2.1	Starting a new thread	24
2.2	Starting a new thread in daemon mode	25
2.3	Workers using multithreading	25
2.4	<code>multiprocessing.Process</code> usage	28
2.5	Result of <code>time python multiprocessing-workers.py</code>	29
2.6	Worker using <code>multiprocessing</code>	30
2.7	Worker using <code>concurrent.futures.ThreadPoolExecutor</code>	31
2.8	Time and output of <code>futures-threads-worker</code>	32
2.9	Worker using <code>concurrent.futures.ProcessPoolExecutor</code>	33
2.10	Extract of <code>concurrent.futures.process</code>	33
2.11	Time and output of <code>futures-threads-worker</code>	34
2.12	Extract of <code>concurrent.futures.thread</code>	34
2.13	Workers using <code>futurist.ThreadPoolExecutor</code>	35
2.14	Output of <code>futures-threads-worker</code>	36
2.15	Using <code>check_and_reject</code> to limit queue size	36
2.16	Using <code>futurist.periodics</code>	37
2.17	Output of <code>futurist-periodics.py</code>	38
2.18	Daemon using <i>Cotyledon</i>	39
2.19	Producer/consumer using <i>Cotyledon</i>	41
2.20	Reconfiguring the number of processes with <i>Cotyledon</i>	43
3.1	A blocking socket	54
3.2	A non-blocking socket	55
3.3	Using <code>select.select</code> with sockets	55

LIST OF EXAMPLES

3.4	Hello world asyncio coroutine	57
3.5	Output of <code>asyncio-basic.py</code>	58
3.6	Coroutine awaiting on coroutine	58
3.7	Output of <code>asyncio-coroutines.py</code>	59
3.8	Coroutine using <code>asyncio.sleep</code> and <code>asyncio.gather</code>	59
3.9	Using <code>aiohttp</code>	61
3.10	Using <code>loop.call_later</code>	62
3.11	Using <code>loop.call_later</code> repeatedly	62
3.12	asyncio TCP server	63
3.13	Using <code>nc</code> to connect to the <code>YellEchoServer</code>	67
3.14	asyncio TCP client	67
3.15	asyncio TCP server with statistics	68
3.16	Testing the asyncio TCP echo server	70
4.1	Using <code>first</code>	82
4.2	Using the operator module with <code>itertools.groupby</code>	85
5.1	Pushing a <code>rq</code> job in a queue	90
5.2	Pushing a <code>rq</code> job in a queue	92
5.3	A simple <i>Celery</i> task	95
5.4	A <i>Celery</i> task with retry support	97
5.5	A <i>Celery</i> task chain	98
5.6	A <i>Celery</i> task called with a specific queue	100
6.1	Retrying pattern	112
6.2	Retrying pattern with <code>sleep</code>	112
6.3	Retrying pattern with exponential backoff	113
6.4	Basic retrying with <i>tenacity</i>	115
6.5	Fixed waiting with <i>tenacity</i>	115
6.6	Exponential back-off waiting with <i>tenacity</i>	116
6.7	Combining wait time with <i>tenacity</i>	117
6.8	Specific retry condition with <i>tenacity</i>	118
6.9	Combining retry condition with <i>tenacity</i>	119

LIST OF EXAMPLES

6.10	Using <i>tenacity</i> without decorator	120
7.1	Threads using a lock	122
7.2	Threads using a reentrant lock	124
7.3	Threads using an <code>threading.Event</code> object	125
7.4	Printing cats in parallel	126
7.5	Printing cats in parallel with a lock	128
7.6	Using <i>Fasteners</i> for interprocess locking	130
7.7	Using <i>Fasteners</i> decorator for interprocess locking	130
7.8	Locking with <code>etcd</code>	133
7.9	Locking with <code>etcd</code> using the <code>with</code> statement	133
7.10	Locking service with <code>etcd</code> and <i>Cotyledon</i>	134
7.11	Getting a <i>Tooz</i> coordinator	136
7.12	Getting a lock	137
7.13	Using <i>Tooz</i> lock and the <code>with</code> statement.	138
8.1	Joining a group	141
8.2	Using <i>Tooz</i> capabilities	144
8.3	Using watchers with <i>Tooz</i>	146
8.4	Example output of Example 8.3	148
8.5	Using <i>Tooz</i> hashing	150
8.6	Using <i>Tooz</i> <code>join_partitioned_group</code> method	155
8.7	Using hash ring to spread Web pages fetching	157
8.8	Running only one Web page fetching program	159
8.9	Running two Web page fetching programs	159
9.1	Basic WSGI application	174
9.2	Basic WSGI application with <code>wsgiref.simple_server</code>	174
9.3	Using WSGI	176
9.4	The <code>PUBLISH</code> command	178
9.5	Checking the <code>SUBSCRIBE</code> command	179
9.6	Receiving messages in Python	179
9.7	Output of <code>listen.py</code>	180

LIST OF EXAMPLES

9.8	<i>Flask</i> based streamer application	180
9.9	ETag header	183
9.10	Using If-None-Match header	184
9.11	Flask application example with ETag usage	184
9.12	Flask application example with ETag usage and PUT	187
9.13	Flask application example with ETag usage and PUT	189
9.14	Flask application example using asynchronous job handling	191
9.15	Using <i>Session</i> with <i>requests</i>	195
9.16	Configuring pool size with <i>requests</i>	195
9.17	Using <i>futures</i> with <i>requests</i>	196
9.18	Using <i>requests-futures</i>	197
9.19	Using <i>aiohttp</i>	198
9.20	Program to compare the performances of different <i>requests</i> usage	199
9.21	Output of <i>requests-comparison.py</i>	201
9.22	Streaming with <i>requests</i>	201
9.23	Streaming with <i>aiohttp</i>	202
9.24	Basic <i>Gabbi</i> test file	203
9.25	Basic <i>Gabbi</i> inclusion in a Flask app	204
9.26	<i>Gabbi</i> test file for Flask ETag application	205
9.27	Running <i>Gabbi</i> using <i>gabbi-run</i>	207
10.1	Profile for Heroku	217
10.2	Deploying a Heroku application	218
10.3	Deploying to Ealstic Beanstalk with <i>eb</i>	222
10.4	<i>app.yaml</i> for Google App Engine	225
10.5	<i>main.py</i> for Google App Engine	225
10.6	Deploying a WSGI application on Google Cloud App	225
11.1	Tests using <i>memcached</i> if available	237
11.2	<i>setup-memcached.sh</i>	237
11.3	<i>pifpaf</i> launching <i>PostgreSQL</i>	240
11.4	Tests using <i>memcached</i> if available with <i>pifpaf</i>	242

LIST OF EXAMPLES

11.5	Tests using <i>memcached</i> fixtures	245
11.6	Application using <i>memcached</i> and fixtures	246
11.7	Application using <i>memcached</i> and fixtures, testing all scenarios	247
12.1	A basic caching example	251
12.2	Using <i>cachetools</i>	252
12.3	Using <i>cachetools</i> to cache Web pages	253
12.4	A basic memoization technique	254
12.5	Using <i>functools.lru_cache</i>	255
12.6	Using <i>cachetools</i> for memoization	256
12.7	Connecting to <i>memcached</i>	257
12.8	Handling missing keys in <i>memcached</i>	258
12.9	Fallback with <i>pymemcache</i>	259
12.10	Counting the number of visitor in <i>memcached</i>	261
12.11	Using CAS in <i>memcached</i>	261
13.1	Basic usage of <i>cProfile</i>	268
13.2	Using the <i>cProfile</i> module with a program	269
13.3	Using <i>memory_profiler</i>	280
13.4	Using <i>slice</i> on <i>memoryview</i> objects	282
13.5	A function defined in a function, disassembled	290
13.6	Disassembling a closure	291

About this Book

Version 1.0 released in 2017.

When I released [The Hacker's Guide to Python](#) in 2014, I had no idea that I would be writing a new book so soon. Having worked on [OpenStack](#) for a few more years, I saw how it is easy to struggle with other aspects of Python, even after being on board for a while.

Nowadays, even if computers are super-fast, no server is fast enough to handle millions of request per second, which is a typical workload we want to use them for. Back in the day, when your application was slow, you just had to optimize it or upgrade your hardware – whichever was cheaper. However, in a world where you may already have done both, you need to be able to scale your application horizontally, i.e., you have to make it run on multiple computers in parallel.

That is usually the start of a long journey, filled with concurrency problems and disaster scenarios.

Developers often dismiss Python when they want to write performance enhancing, and distributed applications. They tend to consider the language to be slow and not suited to that task. Sure, Python is not [Erlang](#), but there's also no need to ditch it for [Go](#) because of everyone *saying* it is faster.

I would like to make you aware, dear reader, that a language is never slow. You would not say that English or French is slow, right? The same applies for programming languages. The only thing that can be slow is the implementation of the language – in Python's case, its reference implementation is CPython.

Indeed CPython can be quite sluggish, and it has its share of problems. Every

implementation of a programming language has its downside. However, I think that the ecosystem of Python can make up for that defect.

Python and everything that evolves around it offer a large set of possibilities to extend your application, so it can manage thousands of requests simultaneously, compensating for its lack of distributed design or, sometimes its "slowness".

Moreover, if you need proof, you can ask companies such as Dropbox, PayPal or Google as they all use Python on a large scale. Instagram has 400 million active users every day and **their whole stack is served using Python and Django**.

In this book, we will discuss how one can push Python further and build applications that can scale horizontally, perform well and remain fast while being distributed. I hope it makes you more productive at Python and allows you to write better applications that are also faster!

Most code in this book targets Python 3. Some snippets might work on Python 2 without much change, but there is no guarantee.

CHAPTER 1

Scaling?

“

Scalability is the capability of a system, network, or process to handle a growing amount of work, or its potential to be enlarged to accommodate that growth.

— Wikipedia

”

When we talk about scaling Python, what we mean is making Python application scalable. However, what is scalability?

According to Wikipedia, scalability is “the capability of a system, network, or process to handle a growing amount of work, or its potential to be enlarged to accommodate that growth”. This definition makes scalability difficult to define as an absolute since no definition applies to all applications.

This book concentrates on methods, technologies, and practice that allow one to make applications fast and able to grow in order to handle more jobs – all of

that using the Python programming language and its major implementation, named CPython.

We are all aware that processors are not becoming faster and faster at a rate where a single threaded application could, one day, be *fast enough* to handle any size workload. That means you need to think about using more than just one processor. Building scalable applications implies that you distribute the workload across multiple workers using multiple processing units.

Dividing up the tasks at hand, those workers run across several processors, and in some cases, across several computers.

That is a *distributed application*.

There are fundamental properties to understand about distributed systems before digging into how to build them in Python – or any other language.

We can lay out the following options when writing an application:

- Write a single-threaded application. This should be your first pick, and indeed it implies no distribution. They are the simplest of all applications. They are easy to understand and therefore easier to maintain. However, they are limited by the power of using a single processor.
- Write a multi-threaded application. Most computers – even your smartphone – are now equipped with multiple processing units. If an application can overload an entire CPUs, it needs to spread its workload over other processors by spawning new threads (or new processes). Multi-threading applications are more error-prone than single-threaded applications, but they offer fewer failure scenarios than multi-nodes applications, as no network is involved.
- Write network distributed applications. This is your last resort when your application needs to scale significantly, and not even one big computer with plenty of CPUs is enough. Those are the most complicated applications to write as they use a network. It means they should handle a lot of scenarios, such as a total or partial

1.1. ACROSS CPUS

failure of a node or the network, high latency, messages being lost, and any other terrible property related to the unreliability of networks.

The properties of distribution vary widely depending on the type you pick. Operations on a single processor can be regarded as fast, with low latency while being reliable, and ordered, whereas operations across several nodes should be considered, slow, with high latency. They are often unreliable and unordered.

Consider each architecture choice or change carefully. As seen throughout this book, there are various tools and methods in Python available for dealing with any of those choices. They help to build distributed systems, and therefore scalable applications.

1.1 Across CPUs

Scaling across processors is usually done using multithreading. Multithreading is the ability to run code in parallel using *threads*. Threads are usually provided by the operating system and are contained in a single process. The operating system is responsible to schedule their execution.

Since they run in parallel, that means they can be executed on separate processors even if they are contained in a single process. However, if only one CPU is available, the code is split up and run sequentially.

Therefore, when writing a multithreaded application, the code always runs concurrently but runs in parallel only if there is more than one CPU available.

This means that multithreading looks like a good way to scale and parallelize your application on one computer. When you want to spread the workload, you start a new thread for each new request instead of handling them one at a time.

However, this does have several drawbacks in Python. If you have been in the Python world for a long time, you have probably encountered the word *GIL*, and

know how hated it is. The GIL is the Python *global interpreter lock*, a lock that must be acquired each time *CPython* needs to execute byte-code. Unfortunately, this means that if you try to scale your application by making it run multiple threads, this global lock always limits the performance of your code, as there are many conflicting demands. All your threads try to grab it as soon as they need to execute Python instructions.

The reason that the *GIL* is required in the first place is that it makes sure that some basic Python objects are thread-safe. For example, the code in Example 1.1 would not be thread-safe without the global Python lock.

Example 1.1 Thread-unsafe code without the GIL

```

1 import threading
2
3 x = []
4
5 def append_two(l):
6     l.append(2)
7
8 threading.Thread(target=append_two, args=(x,)).start()
9
10 x.append(1)
11 print(x)

```

That code prints either `[2, 1]` or `[1, 2]`, no matter what. While there is no way to know which thread appends 1 or 2 before the other, there is an assumption built into Python that each `list.append` operation is **atomic**. If it was not atomic, a memory corruption might arise and the list could simply contain `[1]` or `[2]`.

This phenomenon happens because only one thread is allowed to execute a **bytecode** instruction at a time. That also means that if your threads run a lot of bytecodes, there are many contentions to acquire the GIL, and therefore your program cannot be faster than a single-threaded version – or it could even be slower.

The easiest way to know if an operation is thread-safe is to know if it translates

1.1. ACROSS CPUS

to a single bytecode instruction¹ or if it uses a basic type whose operations are atomic².

So while using threads seems like an ideal solution at first glance, most applications I have seen running using multiple threads struggle to attain 150% CPU usage – that is to say, 1.5 cores used. With computing nodes nowadays usually not having less than four or eight cores, it is a shame. Blame the GIL.

There is currently an effort underway (named [gilectomy](#)) to remove the GIL in *CPython*. Whether this effort will pay off is still unknown, but it is exciting to follow and see how far it will go.

However, *CPython* is just one – although the most common – of the available Python implementations. [Jython](#), for example, [doesn't have a global interpreter lock](#), which means that it can run multiple threads in parallel efficiently. Unfortunately, these projects by their very natures lag behind *CPython*, and so they are not useful targets.

Multithreading involves several traps, and one of them is that all the pieces of code running concurrently are sharing the same global environment and variables. Reading or writing global variables should be done exclusive by using techniques such as locking, which complicates your code; moreover, it is an infinite source of human errors.

Getting multi-threaded applications right is hard. The level of complexity means that it is a large source of bugs – and considering the little to be gained in general, it is better not to waste too much effort on it.

So are we back to our initial use cases, with no real solutions on offer? Not true – there's another solution you can use: using multiple processes. Doing this is going to be more efficient and easier as we will see in [Chapter 2](#). It is also a first step before spreading across a network.

¹Details about disassembling code and bytecode instruction are provided in [Section 13.4](#).

²The list is provided in the [Python FAQ](#).

1.2 Distributed Systems

“

A distributed system is one in which the failure of a computer you didn't even know existed can render your own computer unusable.

— Lamport (1987)

”

When an application uses all the CPU power of a node, and you cannot add more processors to your server or switch to a bigger server, you need a plan B.

The next step usually involves multiple servers, linked together via a network of some sort. That means the application starts to be *distributed*: running not only on one node but on multiple, connected, nodes. Spreading the workload over different hosts introduces several advantages, such as:

- Horizontal scalability, the ability to add more nodes as more traffic comes in
- Fault tolerance, as if a node goes down, another one can pick up the traffic of the dysfunctioning one

While this sounds awesome, it also introduces major drawbacks:

- As with multithreading, concurrency and parallelism come into play and complicate the workflow (e.g., locking usage)
- What can fail will fail, such as a random node in the middle of an operation or a laggy network, so tolerance for failure must be built-in

All of this means that an application which is going the distributed route expand its complexity while potentially increasing its throughput. Making this kind of architectural decision requires great wisdom.

Python does not offer so many tools for building a distributed system, but its ecosystem has a few good options as seen throughout this book. For example, it can be pretty easy to distribute jobs across several nodes as covered in Chapter 5. Bigger problems such as coordination and synchronization with sibling nodes also have a few solutions, as discussed in Chapter 7 and Chapter 8.

Finally, a great approach to writing distributed systems is to make them purely functional, i.e., without any shared state. That means such applications should not have even a single shared, global variable, across all of its distributed processes. Stateless systems are the easiest ones to distribute and scale, and therefore systems should be designed as such when possible. Chapter 4 talks about functional programming and the mindset behind writing such programs.

1.3 Service-Oriented Architecture

If you've never heard of it, service-oriented architecture is an architectural style where a software design is made up of several independent components communicating over a network. Each service is a discrete unit of functionality that can work autonomously. That means that the problem should be divided up into interacting logical pieces. If we refer back to the different application styles defined in Chapter 1, SOA refers to network distributed applications.

This kind of architecture is not a perfect or magical solution. It has many drawbacks, but it also has many advantages that make it valuable... and so popular these days for building distributed applications.

The service-oriented architecture is not a first-class citizen in Python, though it makes it easy to use and implement – the language being generic and the ecosystem

1.3. SERVICE-ORIENTED ARCHITECTURE

rich enough.

Services built for this kind of architecture should follow a few principles³ among them being **stateless**. That means services must either modify and return the requested value (or an error) while separating their functioning from the state of the data. This is an essential property, as it makes it easier to scale the services horizontally.

Statelessness is a property that is also shared with the functional programming paradigm, as discussed in Chapter 4. Both of these are relevant topics and principles to know about when designing scalable applications.

How to split your application into different services might deserve a book on its own, but there are mainly two categories:

- Object-oriented approach: each *noun* is a service, e.g., catalog service, phone service, queue service, etc. Such service types are a good way to represent data types.
- Functional approach: each *verb* is a service, e.g., search service, authentication service, crawl service, etc. Such service types are a good way to represent transformations.

Having too many services has a cost, as they come with some overhead. Think of all of the costs associated, such as maintenance and deployment, and not only development time. Splitting an application should always be a well-thought out decision.

If you know that some services need to scale independently, you should probably split them. However, if they are latency sensitive and should work together very closely, involving a lot of communication, that might be where the line is drawn.

Software production is both a technical and social artifact: there might be some services that come naturally to mind due to the social organization of your project.

³Wikipedia offers [a great list of those principles](#).

1.3. SERVICE-ORIENTED ARCHITECTURE

For example, some teams might be responsible for the user database, so that might be their job to create and maintain an independent user service that other components can use to get information and authenticate it. This is also important to take into consideration when choosing where to set the boundaries of your different services.

Once this is all set ⁴, the technical aspect of the implementation comes to mind. Nowadays, the most common type of services that are encountered are Web services, base on the well-known and ubiquitous HTTP protocol. This is what will be largely discussed in Chapter 9.

⁴Of course no architecture is written in stone, and everything can evolve, likewise social groups change and services might come and go.

CHAPTER 2



CPU Scaling

As CPUs are not getting infinitely faster, using multiple CPUs is the best path towards scalability. That means introducing concurrency and parallelism into your program, and that is not an easy task. However, once correctly done, it really does increase the total throughput.

Python offers two options to spread your workload across multiple local CPUs: threads or processes. They both come with challenges; some are not specifically tied to Python, while some are only relevant to its main implementation, i.e., *CPython*.

2.1 Using Threads

Threads in Python are a good way to run a function concurrently other functions. If your system does not support multiple processors, the threads will be executed one after another as scheduled by the operating system. However, if multiple CPUs are available, threads could be scheduled on multiple processing units, once again as determined by the operating system.

By default, there is only one thread – the main thread – and it is the thread

2.1. USING THREADS

that runs your Python application. To start another thread, Python provides the `threading` module.

Example 2.1 Starting a new thread

```
1 import threading
2
3
4 def print_something(something):
5     print(something)
6
7
8 t = threading.Thread(target=print_something, args=("hello",))
9 t.start()
10 print("thread started")
11 t.join()
```

If you run the program in Example 2.1 multiple times, you will notice that the output might be different each time. On my laptop, doing this gives the following:

```
1 $ python examples/chapter2-cpu-scaling/threading-start.py
2 hellothread started
3 $ python examples/chapter2-cpu-scaling/threading-start.py
4 hello
5   thread started
6 $ python examples/chapter2-cpu-scaling/threading-start.py
7 hello
8   thread started
```

If you specifically expected any one of the outputs each time, then you forgot that there is no guarantee regarding the order of execution for the threads.

Once started, the threads `join`: the main thread waits for the second thread to complete by calling its `join` method. Using `join` is handy in terms of not leaving any threads behind.

If you do not `join` all your threads and wait for them to finish, it is possible that the main thread finishes and exits before the other threads. If this happens, your

2.1. USING THREADS

program will appear to be blocked and will not respond to even a simple KeyboardInterrupt signal.

To avoid this, and because your program might not be in a position to wait for the threads, you can configure threads as *daemons*. When a thread is a daemon, it is considered as a background thread by Python and is terminated as soon as the main thread exists.

Example 2.2 Starting a new thread in daemon mode

```
1 import threading
2
3
4 def print_something(something):
5     print(something)
6
7
8 t = threading.Thread(target=print_something, args=("hello",))
9 t.daemon = True
10 t.start()
11 print("thread started")
```

In Example 2.2 , there is no longer a need to use the join method since the thread is set to be a daemon.

The program below is a simple example, which sums one million random integers eight times, spread across eight threads at the same time.

Example 2.3 Workers using multithreading

```
1 import random
2 import threading
3
4 results = []
5
6
7 def compute():
```

2.1. USING THREADS

```
8     results.append(sum(
9         [random.randint(1, 100) for i in range(1000000)]))
10
11
12 workers = [threading.Thread(target=compute) for x in range(8)]
13 for worker in workers:
14     worker.start()
15 for worker in workers:
16     worker.join()
17 print("Results: %s" % results)
```

Running Example 2.3 program returns the following:

```
1 $ time python multithreading-worker.py
2 Results: [50505811, 50471217, 50531481, 50460206, 50462903, 50533718, ←
   50500182, 50480848]
3 python examples/multithreading-worker.py 19.84s user 6.32s system 116% ←
   cpu 22.501 total
```

The program ran on an idle quad cores CPU, which means that Python could have used up to 400% CPU power. However, it was unable to do that, even with eight threads running in parallel – it stuck at 116%, which is just 29% of the hardware’s capabilities.

The graph in Figure 2.1 illustrates that bottleneck: to access all of the system’s CPU, you need to go through CPython’s GIL.

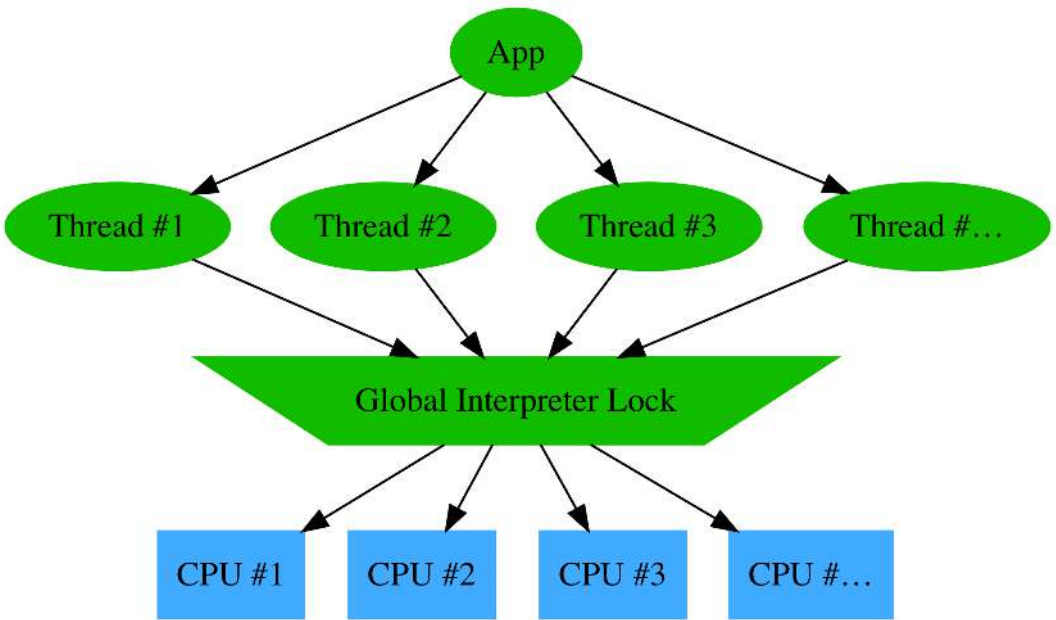


Figure 2.1: Using threads with CPython

Again, as discussed in Section 1.1, the GIL limits the performance of *CPython* when executing multiple threads. Threads are therefore useful when doing parallel computing or input/output on slow networks or files: those tasks can run in parallel without blocking the main thread.

To achieve a greater throughput using multiple CPUs, using processes is an interesting alternative discussed in Section 2.2.

2.2 Using Processes

Since multithreading is not a perfect scalability solution because of the *GIL*, using processes instead of threads is a good alternative. Python obviously exposes the `os`.

2.2. USING PROCESSES

fork system call to create new processes. However, this approach is a little bit too low-level to be interesting in most cases.

Instead, the **multiprocessing** package is a good higher-level alternative. It provides an interface that starts new processes, whatever your operating system might be.

We can rewrite Example 2.3 using processes thanks to the *multiprocessing* library, as shown in Example 2.4.

Example 2.4 multiprocessing.Process usage

```
1 import random
2 import multiprocessing
3
4
5 def compute(results):
6     results.append(sum(
7         [random.randint(1, 100) for i in range(1000000)]))
8
9
10 with multiprocessing.Manager() as manager:
11     results = manager.list()
12     workers = [multiprocessing.Process(target=compute, args=(results,) ←
13         )
14                 for x in range(8)]
15     for worker in workers:
16         worker.start()
17     for worker in workers:
18         worker.join()
19     print("Results: %s" % results)
```

The example is a bit trickier to write as there is no data shared available between different processes. Since each process is a new independent Python, the data is **copied** and each process has its own independent global state. The `multiprocessing.Manager` class provides a way to create shared data structures that are safe for concurrent accesses.

2.2. USING PROCESSES

Running this program gives the following result:

Example 2.5 Result of `time python multiprocessing-workers.py`

```
1 $ time python multiprocessing-workers.py
2 Results: [50505465, 50524237, 50492168, 50482321, 50503634, 50543646, ←
   50533775, 50521610]
3 python examples/multiprocessing-workers.py 32.00s user 0.50s system ←
   332% cpu 9.764 total
```

Compared to Example 2.3, using multiple processes reduces the execution times by 60%. This time, the processes have been able to consume up to 332% of the CPU power, which is more than 80% of the computer's CPU capacity, or close to three times more than multithreading.

The graph in Figure 2.2 tries to lay out the differences in terms of how scheduling processes work and why it is more efficient than using threads, as shown previously.

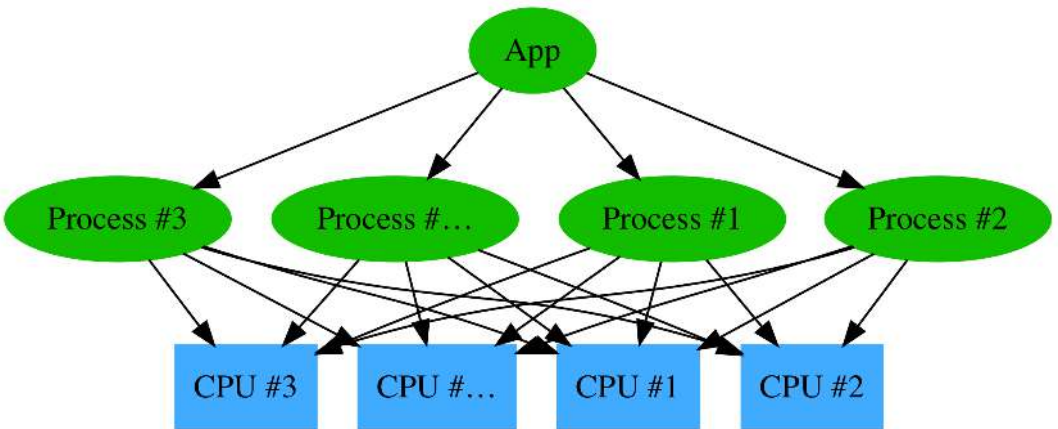


Figure 2.2: Using processes with CPython

Each time some work can be **parallelized** for a certain amount of time, it's

2.3. USING FUTURES

much better to rely on multiprocessing and to fork jobs, thus spreading the workload among several CPU cores, rather than using the threading module.

The multiprocessing library also provides a *pool* mechanism that is useful to rewrite the code from Example 2.4 in a more *functional* manner; an example is provided as Example 2.6.

Example 2.6 Worker using multiprocessing

```
1 import multiprocessing
2 import random
3
4
5 def compute(n):
6     return sum(
7         [random.randint(1, 100) for i in range(1000000)]
8     )
9
10 # Start 8 workers
11 pool = multiprocessing.Pool(processes=8)
12 print("Results: %s" % pool.map(compute, range(8)))
```

Using `multiprocessing.Pool`, there is no need to manage the processes “manually”. The pool starts processes on-demand and takes care of reaping them when done. They are also reusable, which avoids calling the `fork` syscall too often – which is quite costly. It is a convenient design pattern that is also leveraged in futures, as discussed in Section 2.3.

2.3 Using Futures

Python 3.2 introduced the `concurrent.futures` module, which provides an easy way to schedule asynchronous tasks. The module is also available in Python 2 as it has been back-ported – it can easily be installed by running `pip install futures`.

2.3. USING FUTURES

The `concurrent.futures` module is pretty straightforward to use. First, one needs to pick an *executor*. An executor is responsible for scheduling and running asynchronous tasks. It can be seen as a type of engine for execution. The module currently provides two kinds of executors: `concurrent.futures.ThreadPoolExecutor` and `concurrent.futures.ProcessPoolExecutor`. As one might guess, the first one is based on threads and the second one on processes.

As outlined in Section 1.1, the process based executor is going to be much more efficient for long-running tasks that benefit from having an entire CPU available. The threading executor suffers from the same limitation of the threading module, which was covered earlier.

So what is interesting with the `concurrent.futures` module is that it provides an easier to use abstraction layer on top of the `threading` and `multiprocessing` modules. It allows one to run and parallelize code in a straightforward way, providing an abstract data structure called a `concurrent.futures.Future` object.

Each time a program schedules some tasks to execute in threads or processes, the `concurrent.futures` module returns a `Future` object for each of the task scheduled. This `Future` object owns the promise of the work to be completed. Once that work is achieved, the result is available in that `Future` object – so in the end, it does represent the future and the promise of a task to be performed. That is why it is called `Future` in Python, and sometimes *promise* in other languages.

Example 2.7 Worker using `concurrent.futures.ThreadPoolExecutor`

```
1 from concurrent import futures
2 import random
3
4
5 def compute():
6     return sum(
7         [random.randint(1, 100) for i in range(1000000)])
8
9
```

2.3. USING FUTURES

```
10 with futures.ThreadPoolExecutor(max_workers=8) as executor:
11     futures = [executor.submit(compute) for _ in range(8)]
12
13 results = [f.result() for f in futures]
14
15 print("Results: %s" % results)
```

Compared to the threading based example script, you might notice that this one is more functional. I changed the `compute` function to *return* the result rather than changing a shared object. It is then easy to manipulate and transfer the `Future` object and collect the result as desired when it is needed. Functional programming is a perfect paradigm to embrace when trying to spread workload across distributed workers – it is covered more in Chapter 4.

The code just schedules the jobs to be fulfilled and collects the results from the `Future` objects using the `result` method – which also supports a `timeout` parameter in case the program cannot hang for too long. `Future` objects offer some more interesting methods:

- `done()`: This returns `True` if the call was successfully canceled or terminated correctly.
- `add_done_callback(fn)`: This attaches a callable to the future which is called with the future as its only argument; it is done as soon as the future is canceled or terminates correctly.

Example 2.8 Time and output of futures-threads-worker

```
1 $ time python futures-threads-worker.py
2 Results: [50532744, 50524277, 50507195, 50501211, 50537292, 50490570, ←
   50484569, 50515144]
3 python futures-threads-worker.py 14.50s user 6.91s system 126% cpu ←
   16.893 total
```

2.3. USING FUTURES

The execution time is in the same low range as the example using the threading technique: indeed, the underlying engine is based on the threading module.

Keep in mind that `concurrent.futures` allows you to easily switch from threads to processes by using the `concurrent.futures.ProcessPoolExecutor`:

Example 2.9 Worker using `concurrent.futures.ProcessPoolExecutor`

```
1 from concurrent import futures
2 import random
3
4
5 def compute():
6     return sum(
7         [random.randint(1, 100) for i in range(1000000)])
8
9
10 with futures.ProcessPoolExecutor() as executor:
11     futures = [executor.submit(compute) for _ in range(8)]
12
13 results = [f.result() for f in futures]
14
15 print("Results: %s" % results)
```

There is no need to set the number of `max_workers`: as by default `concurrent.futures` calls the `multiprocessing.cpu_count` function to set the number of workers to use, which is equal to the number of CPUs the system can use – as is shown in [Example 2.10](#)

Example 2.10 Extract of `concurrent.futures.process`

```
1 class ProcessPoolExecutor(_base.Executor):
2     def __init__(self, max_workers=None):
3         # [...]
4         if max_workers is None:
5             self._max_workers = multiprocessing.cpu_count()
6         else:
```

2.3. USING FUTURES

```
self._max_workers = max_workers
```

As expected, using processes is much faster than the threading based executor:

Example 2.11 Time and output of futures-threads-worker

```
1 $ time python futures-processes-worker.py
2 Results: [50485099, 50461662, 50553224, 50458097, 50520276, 50510314, ←
   50510035, 50525335]
3 python futures-processes-worker.py 19.48s user 0.30s system 330% cpu ←
   5.991 total
```

Warning

One important thing to notice with both of the pool based executors is the way they manage the processes and threads they spawn.



There are several policies that the authors could have implemented. The one selected is that for each job submitted, a new worker is spawned to do the work, and the work is put in a queue shared across all the existing workers. That means that if the caller sets `max_workers` to 20, then 20 workers will exist as soon as 20 jobs are submitted. None of those processes will ever be destroyed. This is different than, for example, Apache httpd workers that exit after being idle for a while. You can see that this is marked as a TODO in Python source code as shown in [Example 2.12](#)

Example 2.12 Extract of `concurrent.futures.thread`

```
1 class ThreadPoolExecutor(_base.Executor):
2     def submit(self, fn, *args, **kwargs):
3         [...]
4         self._adjust_thread_count()
5
```

```

6 def _adjust_thread_count(self):
7     [...]
8     # TODO(bquinlan): Should avoid creating new threads if there ←
9     # are more
10    # idle threads than items in the work queue.
11    if len(self._threads) < self._max_workers:
12        t = threading.Thread(target=_worker,
13                             args=(weakref.ref(self, weakref_cb),
14                                   self._work_queue))
14        t.daemon = True
15        t.start()
16        self._threads.add(t)
17        _threads_queues[t] = self._work_queue

```

2.4 Advanced Futures Usage

As we have seen in Section 2.3, Future objects are an easy way to parallelize tasks in your application. The *futurist* library has been built on top of *concurrent.futures* and offers a few bonuses that I would like to introduce here. It is (almost) a transparent replacement for *concurrent.futures*, so any code should be straightforward in terms of adapting to this library, which is itself entirely based on *concurrent.futures*.

Example 2.13 Workers using `futurist.ThreadPoolExecutor`

```

1 import futurist
2 from futurist import waiters
3 import random
4
5
6 def compute():
7     return sum(
8         [random.randint(1, 100) for i in range(10000)])
9
10

```

2.4. ADVANCED FUTURES USAGE

```
11 with futurist.ThreadPoolExecutor(max_workers=8) as executor:
12     futures = [executor.submit(compute) for _ in range(8)]
13     print(executor.statistics)
14
15 results = waiters.wait_for_all(futures)
16 print(executor.statistics)
17
18 print("Results: %s" % [r.result() for r in results.done])
```

Example 2.14 Output of futures-threads-worker

```
1 $ python examples/futurist-threads-worker.py
2 <ExecutorStatistics object at 0x10b95b820 (failures=0, executed=0, ←
   runtime=0.00, cancelled=0)>
3 <ExecutorStatistics object at 0x10b95b820 (failures=0, executed=8, ←
   runtime=143.76, cancelled=0)>
4 Results: [50458683, 50479504, 50517520, 50510116, 50450298, 50510857, ←
   50530137, 50511422]
```

First, *futurist* allows any application to access statistics about the executor it uses. That view is valuable for tracking the current status of your tasks and to report information on how the code runs.

futurist also allows passing a function and possibly denying any new job to be submitted by using the `check_and_reject` argument. This argument allows controlling the maximum size of the queue in order to avoid any memory overflow.

Example 2.15 Using `check_and_reject` to limit queue size

```
1 import futurist
2 from futurist import rejection
3 import random
4
5
6 def compute():
7     return sum(
8         [random.randint(1, 100) for i in range(1000000)])
```

2.4. ADVANCED FUTURES USAGE

```
9
10
11 with futurist.ThreadPoolExecutor(
12     max_workers=8,
13     check_and_reject=rejection.reject_when_reached(2)) as executor:
14     futures = [executor.submit(compute) for _ in range(20)]
15     print(executor.statistics)
16
17 results = [f.result() for f in futures]
18 print(executor.statistics)
19
20 print("Results: %s" % results)
```

Depending on the speed of your computer, it is likely that Example 2.15 raises a `futurist.RejectedSubmission` exception because the executor is not fast enough to absorb the backlog, the size of which is limited to two. This example does not catch the exception – obviously, any decent program should handle that exception and either retry later, or raise a different exception to the caller.

futurist addresses a widespread use case with the `futurist.periodics.PeriodicWorker` class. It allows scheduling functions to run regularly, based on the system clock.

Example 2.16 Using `futurist.periodics`

```
1 import time
2
3 from futurist import periodics
4
5
6 @periodics.periodic(1)
7 def every_one(started_at):
8     print("1: %s" % (time.time() - started_at))
9
10
11 w = periodics.PeriodicWorker([
12     (every_one, (time.time(),), {}),
```

2.4. ADVANCED FUTURES USAGE

```
13 ])  
14  
15  
16 @periodics.periodic(4)  
17 def print_stats():  
18     print("stats: %s" % list(w.iter_watchers()))  
19  
20  
21 w.add(print_stats)  
22 w.start()
```

Example 2.17 Output of futurist-periodics.py

```
1 $ python examples/futurist-periodics.py  
2 1: 1.00364780426  
3 1: 2.00827693939  
4 1: 3.00964093208  
5 stats: [<Watcher object at 0x1104fc790 (runs=3, successes=3, failures ←  
    =0, elapsed=0.00, elapsed_waiting=0.00)>,  
6     <Watcher object at 0x1104fc810 (runs=0, successes=0, failures ←  
    =0, elapsed=0.00, elapsed_waiting=0.00)>]  
7 1: 4.00993490219  
8 1: 5.01245594025  
9 1: 6.01481294632  
10 1: 7.0150718689  
11 stats: [<Watcher object at 0x1104fc790 (runs=7, successes=7, failures ←  
    =0, elapsed=0.00, elapsed_waiting=0.00)>,  
12     <Watcher object at 0x1104fc810 (runs=1, successes=1, failures ←  
    =0, elapsed=0.00, elapsed_waiting=0.00)>]  
13 1: 8.01587891579  
14 1: 9.02099585533  
15 [...]
```


2.5. DAEMON PROCESSES

Example 2.16 implements two tasks. One runs every second and prints the time elapsed since the start of the task. The second task runs every four seconds and prints statistics about the running of those tasks. Again here, *futurist* offers internal access to its statistics, which is very handy for reporting the status of the application.

While not necessary to depend on, *futurist* is a great improvement over `concurrent.futures` if you need fine grained control over the execution of your threads or processes.

2.5 Daemon Processes

Being aware of the difference between multithreading and multiprocessing in Python, it becomes more clear that using multiple processes to schedule different jobs is efficient. A widespread use case is to run long-running, background processes (often called daemons) that are responsible for scheduling some tasks regularly or processing jobs from a queue.

It could be possible to leverage `concurrent.futures` and a `ProcessPoolExecutor` to do that as discussed in Section 2.3. However, the pool does not provide any control regarding how it dispatches jobs. The same goes for using the `multiprocessing` module. They both make it hard to efficiently control the running of background tasks. Think of it as the "pets vs. cattle" analogy for processes.

In this section, I would like to introduce you to *Cotyledon*, a Python library designed to build long-running processes.

Example 2.18 Daemon using *Cotyledon*

```
1 import threading
2 import time
3
4 import cotyledon
5
```

2.5. DAEMON PROCESSES

```
6
7 class PrinterService(cotyledon.Service):
8     name = "printer"
9
10    def __init__(self, worker_id):
11        super(PrinterService, self).__init__(worker_id)
12        self._shutdown = threading.Event()
13
14    def run(self):
15        while not self._shutdown.is_set():
16            print("Doing stuff")
17            time.sleep(1)
18
19    def terminate(self):
20        self._shutdown.set()
21
22
23 # Create a manager
24 manager = cotyledon.ServiceManager()
25 # Add 2 PrinterService to run
26 manager.add(PrinterService, 2)
27 # Run all of that
28 manager.run()
```

Example 2.18 is a simple implementation of a daemon using *Cotyledon*. It creates a class named `PrinterService` that implements the needed method for `cotyledon.Service`: `run` which contains the main loop, and `terminate`, which is called by another thread when it terminates the service.

Cotyledon uses several threads internally (at least to handle signals), which is why the `threading.Event` object is used to synchronize the `run` and `terminate` methods.

This service does not do much; it simply prints the message `Doing stuff` every second. The service is started twice by passing two as the number of services to start to `manager.add`. That means *Cotyledon* starts two processes, each of them launching the `PrinterService.run` method.

2.5. DAEMON PROCESSES

When launching this program, you can run the `ps` command on your system – on Unix at least – to see what is running:

```
1 74476 ttys004 0:00.09 cotyledon-simple.py: master process [examples/ ←  
cotyledon-simple.py]  
2 74478 ttys004 0:00.00 cotyledon-simple.py: printer worker(0)  
3 74479 ttys004 0:00.00 cotyledon-simple.py: printer worker(1)
```

Cotyledon runs a master process that is responsible for handling all of its children. It then starts the two instances of `PrinterService` as it was requested to launch. It also gives them nice shiny process names, making them easier to track in the long list of processes. If one of the processes gets killed or crashes, it is automatically relaunched by *Cotyledon*. The library does a lot behind the scenes, e.g., doing the `os.fork` calls and setting up the right modes for daemons.

Cotyledon also supports all operating systems supported by Python itself, avoiding the developer needing to have to think about operating system portability – which can be quite complex.

Example 2.18 is a simple scenario for independent workers – they can execute a job on their own, and they do not need to communicate with each other. This scenario is rare, as most services need to exchange between one another.

Example 2.19 shows an implementation of the common producer/consumer pattern. In this pattern, a service fills a queue (the producer) and other services (the consumers) consume the jobs to execute them.

Example 2.19 Producer/consumer using *Cotyledon*

```
1 import multiprocessing  
2 import time  
3  
4 import cotyledon  
5  
6
```

2.5. DAEMON PROCESSES

```
7 class Manager(cotyledon.ServiceManager):
8     def __init__(self):
9         super(Manager, self).__init__()
10        queue = multiprocessing.Manager().Queue()
11        self.add(ProducerService, args=(queue,))
12        self.add(PrinterService, args=(queue,), workers=2)
13
14
15 class ProducerService(cotyledon.Service):
16     def __init__(self, worker_id, queue):
17         super(ProducerService, self).__init__(worker_id)
18         self.queue = queue
19
20     def run(self):
21         i = 0
22         while True:
23             self.queue.put(i)
24             i += 1
25             time.sleep(1)
26
27
28 class PrinterService(cotyledon.Service):
29     name = "printer"
30
31     def __init__(self, worker_id, queue):
32         super(PrinterService, self).__init__(worker_id)
33         self.queue = queue
34
35     def run(self):
36         while True:
37             job = self.queue.get(block=True)
38             print("I am Worker: %d PID: %d and I print %s"
39                 % (self.worker_id, self.pid, job))
40
41
42 Manager().run()
```

2.5. DAEMON PROCESSES

The program in Example 2.19 implements a custom `cotyledon.ServiceManager` that is in charge for creating the queue object. This queue object is passed to all the services. The `ProducerService` uses that queue and fills it with an incremented integer every second, whereas the `PrinterService` instances consume from that queue and print its content.

When run, the program outputs the following:

```
1 I am Worker: 0 PID: 24727 and I print 0
2 I am Worker: 0 PID: 24727 and I print 1
3 I am Worker: 1 PID: 24728 and I print 2
4 I am Worker: 0 PID: 24727 and I print 3
```

The `multiprocessing.queues.Queue` object eases the communication between different processes. It is safe to use across threads and processes, as it leverages locks internally to guarantee data safety.

Note



If you are familiar with the *Go* programming language, this is the basic pattern that is used to implement the **Go routines** and their channels. That common and efficient pattern made the *Go* language very popular. In *Go*, forking new processes and passing messages between them is provided as a built-in element of the language. Providing syntactic sugar makes it quicker to write programs with this pattern. However, in the end, you can achieve the same thing in *Python*, though, with maybe, a little more effort.

Last, but not least, *Cotyledon* also offers a few more features, such as reloading the program configuration or changing the number of workers for a class dynamically.

Example 2.20 Reconfiguring the number of processes with *Cotyledon*

2.5. DAEMON PROCESSES

```
1 import multiprocessing
2 import time
3
4 import cotyledon
5
6
7 class Manager(cotyledon.ServiceManager):
8     def __init__(self):
9         super(Manager, self).__init__()
10        queue = multiprocessing.Manager().Queue()
11        self.add(ProducerService, args=(queue,))
12        self.printer = self.add(PrinterService, args=(queue,), workers ←
13            =2)
14        self.register_hooks(on_reload=self.reload)
15
16    def reload(self):
17        print("Reloading")
18        self.reconfigure(self.printer, 5)
19
20 class ProducerService(cotyledon.Service):
21     def __init__(self, worker_id, queue):
22         super(ProducerService, self).__init__(worker_id)
23         self.queue = queue
24
25     def run(self):
26         i = 0
27         while True:
28             self.queue.put(i)
29             i += 1
30             time.sleep(1)
31
32
33 class PrinterService(cotyledon.Service):
34     name = "printer"
35
```

2.5. DAEMON PROCESSES

```
36 def __init__(self, worker_id, queue):
37     super(PrinterService, self).__init__(worker_id)
38     self.queue = queue
39
40 def run(self):
41     while True:
42         job = self.queue.get(block=True)
43         print("I am Worker: %d PID: %d and I print %s"
44             % (self.worker_id, self.pid, job))
45
46
47 Manager().run()
```

In Example 2.20, only two processes for `PrinterService` are started. As soon as `SIGHUP` is sent to the master process, `Cotyledon` calls the `Manager.reload` method that reconfigure the printer service to now have five processes. This is easy to check:

```
1 $ ps ax | grep cotyledon
2 55530 s002 S+    0:00.12 cotyledon-reconfigure.py: master process [ ←
   examples/cotyledon-reconfigure.py]
3 55531 s002 S+    0:00.02 cotyledon-reconfigure.py: master process [ ←
   examples/cotyledon-reconfigure.py]
4 55532 s002 S+    0:00.01 cotyledon-reconfigure.py: ProducerService ←
   worker(0)
5 55533 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(0)
6 55534 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(1)
7 $ kill -HUP 55530
8 $ ps ax | grep cotyledon
9 55530 s002 S+    0:00.27 cotyledon-reconfigure.py: master process [ ←
   examples/cotyledon-reconfigure.py]
10 55531 s002 S+    0:00.03 cotyledon-reconfigure.py: master process [ ←
   examples/cotyledon-reconfigure.py]
11 55551 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(2)
12 55553 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(3)
13 55554 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(4)
14 55555 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(1)
15 55557 s002 S+    0:00.01 cotyledon-reconfigure.py: ProducerService ←
```

```
worker(0)
55558 s002 S+    0:00.01 cotyledon-reconfigure.py: printer worker(0)
```

Cotyledon is an excellent library for managing long-running processes. I encourage everyone to leverage it to build long-running, background, job workers.

2.6 Mehdi Abaakouk on CPU Scaling

“



Hey Mehdi! Could you start by introducing yourself and explaining how you came to Python?

Hi! I am Mehdi Abaakouk, I live in Toulouse (France), and I have been using Linux for almost twenty years.

My current job is Senior Software Engineer for Redhat. My main interests in computer sciences are open-source software and how the Internet works under the hood, and I like hacking both of them.

At the beginning of my using Linux, I was frustrated with the music players available at that time, so I started to write one. I looked at the code of many media players and wanted to use GTK/GStreamer toolkits. I first tried it in C by reusing some code from Rhythmbox, but I quickly abandoned that because of the slow progress I made

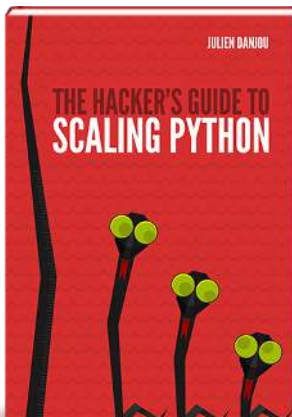
Hey, this was only a sample chapter!

I hope that you did like the sample! It includes the complete table of contents and a full chapter with its examples.

The full version of THE HACKER'S GUIDE TO SCALING PYTHON includes:

- 13 chapters
- 7 interviews
- 300 pages
- 80 code snippets
- Practical examples
- Available in PDF, HTML, EPUB and MOBI formats
- Available in a professionally printed paperback format
- And a few more bonuses such as Docker images!

Buy the Book!



Now that you've read the sample, you might be interested in buying the whole book. It's available online at scaling-python.com in different formats and packages. Go check it out!