# Chapter 7

# Threading

## 7.1 Threading

Threads are the smallest program units that an operating system can execute. Programming with threads allows that several lightweight processes can run simultaneously inside the same program. Threads that are in the same process share the memory and the state of the variables of the process. This shared use of resources enables threads to run faster than execute several instances of the same program.

Each process has at least one thread that corresponds to its execution. When a process creates several threads, it executes these units as *parallel processes*. In a single-core machine, the parallelism is approximated through thread *scheduling* or *time slicing*. The approximation consists of assigning a limited amount of time to each thread repeatedly. The alternation between threads simulates parallelism. Although there is no true increase in execution speed, the program becomes much more responsive. For example, several tasks may execute while a program is waiting for a user input. Multi-core machines achieve a truly faster execution of the program. Figure 7.1 shows how threads interact with the main process.

Some examples of where it is useful to implement threads, even on single-core computers, are:

- Interfaces that interact with the user while the machine executes a heavyweight calculation process.

- Delegation of tasks that follow consumer-producer pattern, *i.e.*, jobs which outputs and inputs are related, but run independently.

- Multi-users applications, in which each thread would be in charge of the requests of each user.

Python 3 handles threads by using the *threading* library. It includes several methods and objects to manipulate threads.
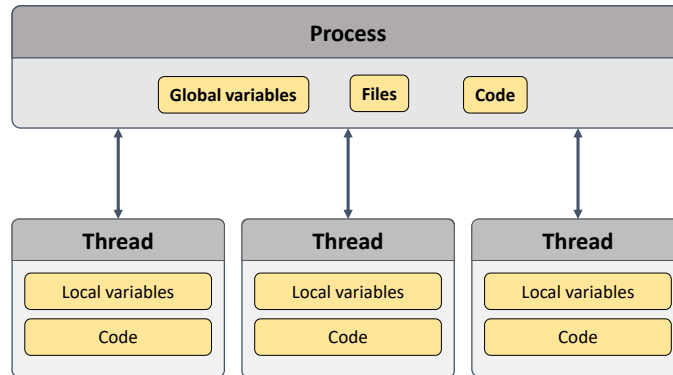
Figure 7.1: Diagram of a threading-based application
.

## Creating Threads

We can create a new thread using the Thread class from the Threading library. This class requires three arguments: target to define the function to be executed; name to provide name we want to give to the thread; args to pass the target arguments. Once created, it may be executed by calling the start() method. In the next example, we create three threads t1, w1, and w2, that execute different instances of the service and worker functions.

```python
# code0.py

import threading
import time


def worker():
    print("{} starting...".format(threading.currentThread().getName()))
    # This stops the thread execution for 2 seconds.
    time.sleep(2)
    print("{} exiting...".format(threading.currentThread().getName()))


def service():
    print("{} starting...".format(threading.currentThread().getName()))
    # This stops the thread execution for 4 seconds.
    time.sleep(4)
    print("{} exiting...".format(threading.currentThread().getName()))
```

```python
19
20
21  # We create two named threads
22  t1 = threading.Thread(name='Thread 1', target=service)
23  w1 = threading.Thread(name='Thread 2', target=worker)
24
25  # This uses the default name (Thread-i)
26  w2 = threading.Thread(target=worker)
27
28  # All threads are executed
29  w1.start()
30  w2.start()
31  t1.start()
32
33
34  # The following will be printed before the threads finish executing
35  print('\nThree threads were created\n')
```

```
Thread 2 starting...
Thread-1 starting...
Thread 1 starting...


Three threads were created


Thread 2 exiting...
Thread-1 exiting...
Thread 1 exiting...
```

In the example, we see that once we have initialized the threads, the main program continues with the rest of the instructions while threads execute their task. The three threads end independently at different times. The main program waits until all the threads finish correctly.

The following code shows an example of how to pass arguments to the `target` function through the `args` attribute.

```python
1  # code1.py
2
3  import threading
```

```python
4  import time
5
6
7  def worker(t):
8      print("{} starting...".format(threading.currentThread().getName()))
9
10     # Thread is stopped for t seconds
11     time.sleep(t)
12     print("{} exiting...".format(threading.currentThread().getName()))
13
14
15 # Threads are created using the Thread class, these are associated with the
16 # objective function to be executed by the thread. Function attributes are
17 # given using the 'args' keyword. In this example, we only need to give one
18 # argument. For this reason a one value tuple is given.
19
20 w = threading.Thread(name='Thread 2', target=worker, args=(3,))
21 w.start()

   Thread 2 starting...
   Thread 2 exiting...
```

Another way of creating a thread is by inheriting from `Thread` and redefining the `run()` method.

```python
1  # code2.py
2
3  import threading
4  import time
5
6
7  class Worker(threading.Thread):
8
9      def __init__(self, t):
10         super().__init__()
11         self.t = t
12
13     def run(self):
```

```python
14          print("{} starting...".format(threading.currentThread().getName()))
15          time.sleep(self.t)
16          print("{} exiting...".format(threading.currentThread().getName()))
17
18
19   class Service(threading.Thread):
20
21       def __init__(self, t):
22           super().__init__()
23           self.t = t
24
25       def run(self):
26           print("{} starting...".format(threading.currentThread().getName()))
27           time.sleep(self.t)
28           print("{} exiting...".format(threading.currentThread().getName()))
29
30
31   # Creating threads
32   t1 = Service(5)
33   w1 = Worker(2)
34   w2 = Worker(4)
35
36   # The created threads are executed
37   t1.start()
38   w1.start()
39   w2.start()
```

```
Thread-1 starting...
Thread-2 starting...
Thread-3 starting...
Thread-2 exiting...
Thread-3 exiting...
Thread-1 exiting...
```

## Join()

In certain situations, we would like to synchronize part of our main program with the outputs of the running threads. When we need the main program to wait that the execution of a thread or a group of threads finished, we must use the `join(< maximum-waiting-time >)` method after the thread starts. In this way, every time we use `join()` the main program will be blocked until the referenced threads finish correctly. If we do not define the maximum waiting time, the main program waits indefinitely until the referenced thread finishes. Figure 7.2 shows the execution of the program using `join()`.
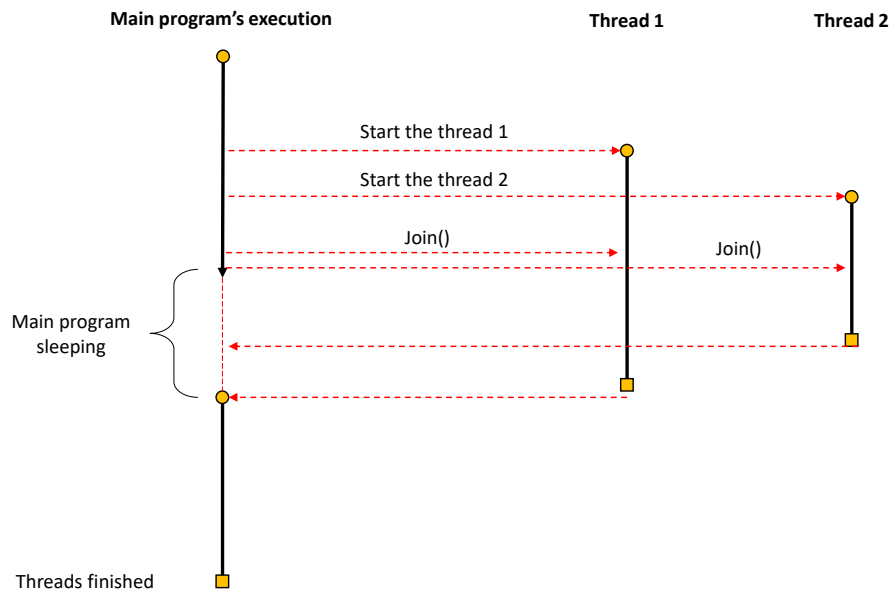


Figure 7.2: Diagram shows the program's flow when we use the `join()` method. We can see that the main program will sleep until *thread 1* finishes. The *thread 2* keeps running independently to the other thread and the main program.

Now let's see the same previous example but incorporating the `join()` method after threads start running.

```
1
2  # Creating threads
3  t1 = Service(5)
4  w1 = Worker(2)
5  w2 = Worker(4)
6
7  # Starting threads
8  t1.start()
9  w1.start()
```

```
10  w2.start()

11

12  # Here we call the join() method to block the main program.

13  # The other threads keep running independently

14  t0 = time.time()

15  w1.join()

16  print('Main program waits for: {}'.format(time.time() - t0))
```

```
Thread 1 starting...
Thread 2 starting...
Thread 3 starting...
Thread 2 exiting...
Main program waits for: 2.000131607055664
Thread 1 exiting...
Thread 3 exiting...
```

### IsAlive()

We can identify if a thread finished its execution using the `IsAlive()` method or the `is_alive` attribute, for example, after using `join()`. The following example shows the way to use `IsAlive()` to check if a thread is still running after a certain amount of time.

```
1

2  t = Service(4)

3  t.start()

4

5  # The main program will wait 5 seconds after 't' has finished executing

6  # before continuing its execution.

7  t.join(5)

8

9  # This returns true if the thread is not currently executing

10  if not t.isAlive():

11      print('The thread has finished successfully')

12  else:

13      print('The thread is still executing')
```

```
Thread-1 starting...
```

```
    Thread-1 exiting...
    The thread has finished successfully
```

We can avoid the use of too many `prints` that help us with the tracking of threads, by using the *logging* library. Every time we make a log we have to embed the name of each thread on its log message, as shown in the following example:

```python
1   # code5.py
2
3   import threading
4   import time
5   import logging
6
7
8   # This sets ups the format in which the messages will be logged on console
9   logging.basicConfig(level=logging.DEBUG, format='[%(levelname)s]'
10                                           '(%(threadName)-10s) %(message)s')
11
12  class Worker(threading.Thread):
13
14      def __init__(self, t):
15          super().__init__()
16          self.t = t
17
18      def run(self):
19          logging.debug('Starting')
20          time.sleep(self.t)
21          logging.debug('Exiting')
22
23
24  class Service(threading.Thread):
25
26      def __init__(self, t):
27          super().__init__()
28          self.t = t
29
```

```python
30      def run(self):
31          logging.debug('Starting')
32          time.sleep(self.t)
33          logging.debug('Exiting')
34
35
36  # Creating threads
37  t1 = Service(4)
38  w1 = Worker(2)
39  w2 = Worker(2)
40
41  # Starting threads
42  w1.start()
43  w2.start()
44  t1.start()
```

```
[DEBUG](Thread-2  ) Starting
[DEBUG](Thread-3  ) Starting
[DEBUG](Thread-1  ) Starting
[DEBUG](Thread-2  ) Exiting
[DEBUG](Thread-3  ) Exiting
[DEBUG](Thread-1  ) Exiting
```

### Daemon Threads

In general, the main program waits for all the threads to finish before ending its execution. *Daemon threads* let the main program to kill them off after other threads (and itself) finish. Daemon threads do not prevent the main program to end. In this way, daemon threads will only run until the main program finishes.

```python
1  # code7.py
2
3  import threading
4  import time
5
6
7  class Worker(threading.Thread):
8
```

```python
 9      def __init__(self, t):
10          super().__init__()
11          self.t = t
12
13      def run(self):
14          print("{} starting...".format(threading.currentThread().getName()))
15          time.sleep(self.t)
16          print("{} exiting...".format(threading.currentThread().getName()))
17
18
19  class Service(threading.Thread):
20
21      def __init__(self, t):
22          super().__init__()
23          self.t = t
24          # We can set a thread as deamon inside the class definition
25          # setting the daemon attribute as True
26          self.daemon = True
27
28      def run(self):
29          print("{} starting...".format(threading.currentThread().getName()))
30          time.sleep(self.t)
31          print("{} exiting...".format(threading.currentThread().getName()))
32
33
34  # Creating threads
35  t1 = Service(5)
36  w1 = Worker(2)
37
38  # Setting the working thread as daemon
39  # We can use this same method when we define a function as target
40  # of a thread.
41  w1.setDaemon(True)
42
43  # Executing threads
```

```
44  w1.start()
45  t1.start()
```

```
    Thread 2 starting...
    Thread 1 starting...
```

The previous example explains the use of daemon threads. The console output shows how threads are interrupted abruptly after the main program ends its execution. We can compare this output with the output of the next example, configuring threads as daemon (removing lines 24 and 39):

```
    Thread-2 starting...
    Thread-1 starting...
    Thread-2 exiting...
    Thread-1 exiting...
```

Note that threads complete the execution and the program did not close until both threads finished. If for any reason, we require waiting for a daemon thread during an amount of time, we can specify that amount (in seconds) in the `join()` method:

```
1   # code9.py
2
3   import threading
4   import logging
5   import time
6
7   logging.basicConfig(level=logging.DEBUG, format='(%(threadName)-10s) '
8                                                    '%(message)s')
9
10  class DaemonThread(threading.Thread):
11
12      def __init__(self, t):
13          super().__init__()
14          self.t = t
15          self.daemon = True
16          self.name = 'daemon'
17
18      def run(self):
```

```python
19          logging.debug('Starting')
20          time.sleep(self.t)
21          logging.debug('Exiting')
22
23  class NonDaemonThread(threading.Thread):
24
25      def __init__(self, t):
26          super().__init__()
27          self.t = t
28          self.name = 'non-daemon'
29
30      def run(self):
31          logging.debug('Starting')
32          time.sleep(self.t)
33          logging.debug('Exiting')
34
35  # Creating threads
36  d = DaemonThread(3)
37  t = NonDaemonThread(1)
38
39  # Executing threads
40  d.start()
41  t.start()
42
43  # Waiting thread d for 1 seconds
44  d.join(2)
45  print('is d alive?: {}'.format(d.isAlive()))
```

```
(daemon    ) Starting
(non-daemon) Starting
(non-daemon) Exiting
is d alive?: True
```

### Timers

The class `Timer` is a subclass of the class `Thread` and allows us to execute a process or an action after a certain amount of time has passed. `Timer` requires as basic parameters the time in seconds after which the thread starts running, the name of the process to execute and the entry arguments for the process. The `cancel()` method allows us, if required, to cancel the execution of the timer before its begins.

```python
1  # code10.py
2
3  import threading
4
5  def delayed_message(msg):
6      print("Message:", msg)
7
8  t1 = threading.Timer(10.0, delayed_message, args=("This is a thread t1!",))
9  t2 = threading.Timer(5.0, delayed_message, args=('This is a thread t2!',))
10 t3 = threading.Timer(15.0, delayed_message, args=('This is a thread t3!',))
11
12 # This thread will start after 10 seconds
13 t1.start()
14
15 # This thread will start after 5 seconds
16 t2.start()
17
18 # Here we cancel thread t1
19 t1.cancel()
20
21 # This thread will start after 15 seconds
22 t3.start()

   Message: This is a thread t2!
   Message: This is a thread t3!
```

## 7.2 Synchronization

Threads run in a non-deterministic way. Therefore, there are some situations in which more than one thread must share the access to certain resources, such as files and memory. During this process, **only one** thread have access to the

resource, and the remaining threads must wait for it. When there is multiple *concurrence* to a resource it is possible to control the access through synchronization mechanisms among the threads.

## Locks

*Locks* allow us to synchronize the access to shared resources between two or more threads. The *Threading* library provides us with the `Lock` class which allows the synchronization. A lock has two states: *locked* and *unlocked*. The default state is *unlocked*. When a given thread $t_i$ attempts to execute, first it tries to acquire the lock (with the `acquire()` method). If another thread $t_j$ takes the lock, $t_i$ must wait for $t_j$ to finish and release it (with the `release()` method) to have the chance to acquire the lock. Once $t_i$ acquires the lock, it can start executing. Figure 7.3 shows a general scheme of synchronization between threads using locks.

```python
1   # code11.py
2
3   import threading
4
5
6   # This class models a thread that blocks to a file
7   class MyThread(threading.Thread):
8
9       lock = threading.Lock()
10
11      def __init__(self, i, file):
12          super().__init__()
13          self.i = i
14          self.file = file
15
16      # This method is the one executed when the start() method is called.
17      def run(self):
18
19          # Blocks other threads from entering the next block
20          MyThread.lock.acquire()
21          try:
22              self.file.write('This line was written by thread #{}\n'.format(self.i))
23          finally:
24              # Releases the resource
```

```python
25                  MyThread.lock.release()
26
27
28  if __name__ == '__main__':
29      n_threads = 15
30      threads = []
31
32      # We create a file to write the output.
33      with open('out.txt', 'w') as file:
34
35          # All writing threads are created at once
36          for i in range(n_threads):
37
38              my_thread = MyThread(i, file)
39
40              # The thread is started, which executes the run() method.
41              my_thread.start()
42              threads.append(my_thread)
```

```
This line was written by thread #0
This line was written by thread #1
This line was written by thread #2
This line was written by thread #3
This line was written by thread #4
This line was written by thread #5
This line was written by thread #6
This line was written by thread #7
This line was written by thread #8
This line was written by thread #9
This line was written by thread #10
This line was written by thread #11
This line was written by thread #12
This line was written by thread #13
This line was written by thread #14
```

Fortunately in Python *locks* can also work inside a *context manager* through the `with` sentence. In this case, is the
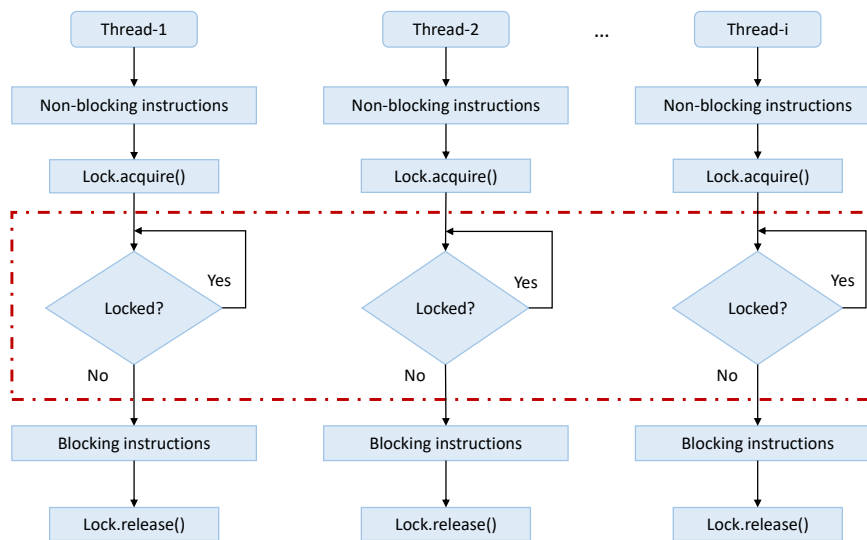
Figure 7.3: This image shows a general example where a set of *i* threads are running. Whenever a thread acquires the lock, the rest of the threads must wait to be able to execute. After the thread releases the lock, the other threads can acquire it and run their (blocking) instructions.

same `with` that is in charge of calling the `acquire()` and `release()` methods. For example, the *locks* used in

the `run` method from the previous example can be implemented as shown:

```
1  def run(self):
2      with MyThread.lock:
3          self.file.write(
4              'This line was written by thread #{}\n'.format(self.i))
```

A common problem in concurrent programming is the *Producer-Consumer* pattern. This problem arises when two or

more threads, known as *producers* and *consumers*, access to the same storage space or *buffer*. Under this scheme,

producers put items in the *buffer* and consumers pull items out of the *buffer*. This model allows the communication

between different threads. In general the *buffer* shared in this model is implemented through a *synchronized queue* or

*secure queue*.

For example, let's assume that we can separate a program that processes a text file with numbers in two independent

threads. The first thread is in charge of reading the file and appending the values to a queue. The second thread stores

into another file the sum of the tuples of numbers previously added into the queue. We communicate both threads

through a synchronized queue implemented as shown next:

```
1   # code13.py
2
3   import collections
4   import threading
5
6   class MyDeque(collections.deque):
7
8       # We inherit from a normal collections module Deque and
9       # we add the locking mechanisms to ensure thread
10      # synchronization
11
12      def __init__(self):
13          super().__init__()
14          # A lock is created for this queue
15          self.lock = threading.Lock()
16
17      def append(self, element):
18
19          # The lock is used within a context manager
20          with self.lock:
21              super().append(element)
22              print('[ADD] queue now has {} elements'.format(len(self)))
23
24      def popleft(self):
25          with self.lock:
26              print('[REMOVE] queue now has {} elements'.format(len(self)))
27              return super().popleft()
```

Now let's see the rest of the implementation of the producer and the consumer. As a recommendation, we encourage the read to try the examples directly in a terminal or using a IDE such as *PyCharm*.

```
1   # code14.py
2
3   import time
4   import threading
5
```

```
 6
 7  class Producer(threading.Thread):
 8      # This thread is implemented as a class
 9
10      def __init__(self, queue):
11          super().__init__()
12          self.queue = queue
13
14      def run(self):
15          # We open the file using a context manager. We explain this in details
16          # in Chapter 10.
17          with open('raw_numbers.txt') as file:
18              for line in file:
19                  values = tuple(map(int, line.strip().split(',')))
20                  self.queue.append(values)
21
22
23  def consumer(queue):
24      # This thread is implemented as a function
25
26      with open('processed_numbers.txt', 'w') as file:
27          while len(queue) > 0:
28              numbers = queue.pop()
29              file.write('{}\n'.format(sum(numbers)))
30
31              # Simulates that the consumer is slower than the producer
32              time.sleep(0.001)
33
34
35  if __name__ == '__main__':
36
37      queue = MyDeque()
38
39      p = Producer(queue)
40      p.start()
```

```
41
42      c = threading.Thread(target=consumer, args=(queue,))
43      c.start()

    [ADD] queue now has 1 elements
    [ADD] queue now has 2 elements
    [ADD] queue now has 3 elements
    [ADD] queue now has 4 elements
    [ADD] queue now has 5 elements
```

## Deadlock

In the context of multithreading-based applications, there is an innocent but dangerous situation in programs that use locks. This case is commonly called *deadlock*. A deadlock occurs when two or more threads are stuck waiting for each other to release a resource. For example, let `FirstProcess` be a thread that acquires a lock `a` and requests for a lock `b` so it can release the lock `a`. Let `SecondProcess` be another thread that already acquired the lock `b` and is waiting for the lock `a` before it releases the lock `b`. We note the deadlock because of our program will be frozen without getting a runtime error or crash. The next code example shows a template of a deadlock according to with the situation described:

```python
1   # code15.py
2
3   import threading
4
5
6   class FirstProcess(threading.Thread):
7
8       def __init__(self, lock_a, lock_b):
9           self.lock_a = lock_a
10          self.lock_b = lock_b
11
12      def run(self):
13          with self.lock_a:
14              # Acquire the first lock
15
16              with self.lock_b:
```

```
17                  # Acquire the second lock for another concurrent task

18

19

20  class SecondProcess(threading.Thread):

21

22      def __init__(self, lock_a, lock_b):

23          self.lock_a = lock_a

24          self.lock_b = lock_b

25

26      def run(self):

27          with self.lock_b:

28                  # Acquire the first lock

29                  # Notice that this thread require the lock_b, that

30                  # could be taken fot other thread previously

31

32              with self.lock_a:

33                      # Acquire the second lock for another concurrent task

34

35

36  lock_a = threading.Lock()

37  lock_b = threading.Lock()

38

39  t1 = FirstProcess(lock_a, lock_b)

40  t2 = SecondProcess(lock_a, lock_b)

41  t1.start()

42  t2.start()
```

We can decrease the risk of a deadlock by restricting the number of locks that a threads can acquire at a time.

## Queue

Fortunately, Python has an optimized library for secure queues management in *producer-consumer* models. The `queue` library has a implemented queue that safely manages multiples concurrences. It is different to the queue implemented in `collections` library used in data structures because that one does not have locks for synchronization.

The main queue methods in the `queue` library are:

- `put()`: Adds an item to the queue (push)

- `get()`: Removes and returns an item from the queue (pop)

- `task_done()`: Requires to be called each time an item has been processed

- `join()`: Blocks the queue until all the items have been processed

Recall the text file processing example shown before. The implementation using the `queue` library is as follows:

```python
# code16.py

import threading
import time
import queue


class Producer(threading.Thread):

    def __init__(self, que):
        super().__init__()
        self.que = que

    def run(self):
        with open('raw_numbers.txt') as file:
            for line in file:
                values = tuple([int(l) for l in line.strip().split(',')])
                self.que.put(values)
                print('[PRODUCER] The queue has {} elements.'.format( \
                    self.que.qsize()))

                # Simulates a slower process
                time.sleep(0.001)


def consumer(que):
    with open('processed_numbers.txt', 'w') as file:
        while True:
```

```
29
30              # A try/except clause is used in order to stop
31              # the consumer once there is no elements left in the
32              # queue. If not for this, the consumer would be executing
33              # for ever
34
35              try:
36
37                  # If no elements are left in the queue, an Empty
38                  # exception is raised
39                  numbers = que.get(False)
40              except queue.Empty:
41                  break
42              else:
43                  file.write('{}\n'.format(sum(numbers)))
44                  que.task_done()
45
46                  # qsize() returns the queue size
47                  print('[CONSUMER] The queue now has {} elements.'.format( \
48                      que.qsize()))
49
50                  # Simulates a complex process. If the consumer was faster
51                  # than the producer, the threads would end abruptly
52                  time.sleep(0.005)
53
54
55  if __name__ == '__main__':
56
57      q = queue.Queue()
58
59      # a producer is created and executed
60      p = Producer(q)
61      p.start()
62
63      # a consumer thread is created and executed with the same queue
```

```
64      c = threading.Thread(target=consumer, args=(q,))
65      c.start()
```

```
[PRODUCER] The queue has 1 elements.
[CONSUMER] The queue now has 0 elements.
[PRODUCER] The queue has 1 elements.
[PRODUCER] The queue has 2 elements.
[PRODUCER] The queue has 3 elements.
[PRODUCER] The queue has 4 elements.
[CONSUMER] The queue now has 3 elements.
[CONSUMER] The queue now has 2 elements.
[CONSUMER] The queue now has 1 elements.
[CONSUMER] The queue now has 0 elements.
```

## 7.3   Hands-On Activities

### Activity 7.1

*ALERT!* Godzilla has arrived at Santiago! Soldiers need to simulate a battle against Godzilla. The simulation will contribute deciding whether it is better to run away or to fight him. With this purpose, Soldiers have given us a report with the specifications that we have to accomplish. These specifications are:

- There is just one Godzilla and several soldiers.

- Each soldier has an attack speed, remaining life (HP), and strength of attack (damage).

- Godzilla attacks every eight seconds, affecting all soldiers by decreasing their HP in three units.

- Each time a soldier attacks Godzilla, it attacks back decreasing one-fourth of the soldiers' attack to his HP.

- The soldiers' attack speed is random between 4 and 19 seconds.

- You must create one new soldier every $x$ seconds, where $x$ has to be previously defined by you.

### Activity 7.2

Congratulations! Thanks to the previous simulation (7.3), the Army has realized of its superiority against Godzilla. Santiago is safe again, or that is what we believe. The truth is that the epic battle has been nothing but a simulation

made by Godzilla to decide if he attacks Santiago or not. Now Santiago will be faced by Mega-Godzilla (Godzilla in its ultimate form), with all of the powers he has not shown before. The task is to simulate the battle between Mega-Godzilla and the Army so it can be written in history books. The simulation must:

- Contain several soldiers and one Mega-Godzilla.

- Each soldier has an attack speed, remaining life (HP), and strength of attack (damage).

- Mega-Godzilla attacks every $N$ seconds, where $N$ is a random number between 3 and 6. $N$ has to be reset after each action.

- Mega-Godzilla attacks in the following ways:

  - Normal attack: Mega-Godzilla stomps affecting all soldiers. This attack causes a three units damage to each soldier.

  - Scaly Skin: Each time a soldier attacks Mega-Godzilla, it attacks back decreasing one-fourth of the soldiers' attack to his HP.

  - Ultimate Mega-Godzilla Super Attack: Mega-Godzilla screams causing a six units damage to every soldier. Also, the scream stuns all the soldiers for 10 seconds. During this period, soldiers can not perform any action.

- Soldiers' attack speed is random between 4 and 19 seconds.

- Soldiers' attack lasts a random number between 1 and 3 seconds.

- Only one soldier can attack Mega-Godzilla at a time.