

Chapter 16

Programming

16.1 Python Programming

16.1.1 Asynchronous Programming

When working with applications that involve heavy I/O operations, such as downloading stock market data, running database queries, these operations may be time intensive and degrade the performance of the application. If we run these operations sequentially, then the performance issues grow linearly in the number of such operations and the application does not scale. If we can exploit concurrency to make these requests together for simultaneous waiting, our application would be much less idle and benefit from less latency. In general, an application's work can be divided into I/O and CPU bound (bound means limiting factor) work. The I/O operations involve a computer's input and output devices, such as keyboard, hard drive and network cards. CPU operations involving the central processing unit, and consists of work such as mathematical computation. Concurrency refers to allowing more than a single task to be handled at the same time. One of the ways this is achieved in Python is through the *asyncio* module, which stands for asynchronous I/O. Asynchronous programming refers to the act of running tasks in the background, separate from the main application. The *asyncio* library utilizes the concurrency model known as a single-threaded event loop. We will explore what these mean later. Additionally, *asyncio* provides utility functions that allow us to work with multiple threads and multiple processes, giving us better control over these workflows in neat abstraction.

Concurrency means that at least two tasks are happening at the same time. This can be rather vague and needs to be clarified. Concurrency does not in fact imply that something is running in parallel. When we say that parallelism exists, it requires that the tasks are both concurrent and executing at the same time. We may be concurrently letting the cake cool in the oven while preparing the table - these two events can be conducted concurrently. When you hire someone to cook the fish while you prepare the jus, the events occur in parallel. Concurrency can be obtained on a single core CPU via task-switching using a mechanism known as preemptive multitasking. Parallelism can only happen on a machine with at least two cores. Parallelism implies concurrency but not the other way around. With multitasking, one may have multiple tasks happening concurrently - except only one of them is executing at one time. Multitasking falls under two types - preemptive multitasking and cooperative multitasking. In the former, the OS (operating system) decides how and when to switch between different tasks via time slicing. In the latter, the running process decides to give up CPU time via explicitly earmarked regions in the code. *asyncio* in particular employs the cooperating multitasking model to achieve concurrency.

As mentioned, it is single threaded (and therefore also single process) and runs a concurrent-but-not-parallel model. Of course, we can combine parallelism with *asyncio*, and we will see how to do that later. Cooperating multitasking has some benefits. One of that is that it is less resource intensive. When the OS switches between running threads or processes, a context switch is said to be involved, and the OS needs to save the state of the running process/thread in order to continue later. This is expensive. By explicitly earmarking the areas for giving up CPU control, the application also has efficiency gains, assuming it is done properly.

So this begs the question of the difference between a process and a thread. A process is an application that has its own memory space that other applications cannot access. Multiple processes can run on a single machine at the same time if there are multiple cores, otherwise they run concurrently via time slicing. The algorithms that determines when the OS triggers task-switching is OS-dependent, and will hence not be discussed in this section, which is intended to be mostly OS-independent. A thread on the other hand, does not have its own memory space. Instead, they share the memory with the same process that created it. A process is always minimally associated with a thread known as the main thread. The other threads created from the main thread is known as worker threads. Threads can run alongside one another on a multi-core CPU, and the OS can also switch between threads via time slicing. One can think of threads as lightweight processes, with a different memory profile. The entry into any Python application begins with a process and its main thread.

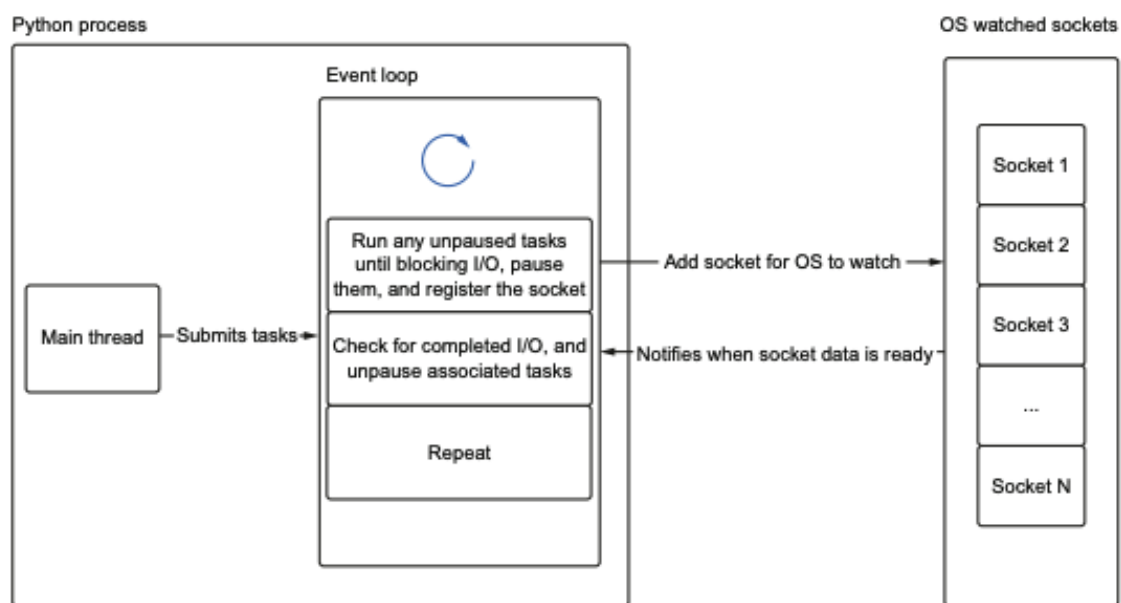
However it turns out that multi-threading for achieving concurrency in Python is somewhat nuanced due to a feature known as the GIL (global interpreter lock). Any high-level language is converted into machine language by compilers or interpreters, or a mix of both, and has lower-level representations in bytecodes, assembly language and so on, which is a further abstraction on top of binary code, which specifies computer instruction by ones and zeros. We are purposefully abstracted away from this gritty detail, and so we will not go into its details. Briefly, the GIL prevents a Python process from executing more than one line of Python bytecode instruction at any time. While this may seem restrictive, the feature exists due to the memory management in CPython, the reference implementation of Python. Memory in CPython is managed by a mechanism known as reference counting, which keeps track the number of variables referencing the object. When a variable pointer is assigned to the object, the reference count is incremented - and then decremented when the variable no longer needs it. When the reference count hits zero, it is put in garbage collection to free up memory. CPython however, is not thread safe, which means that if multiple threads share a variable and attempt to modify it, the result is non-deterministic, a state known as race condition. While threads can run concurrently on different cores, the interpreter lock ensures that only one of them runs Python code at a time to prevent race conditions. The other threads idle. While this can paint a doom picture for Python as a programming language, there are workarounds. Firstly, the GIL is released when the task occurs outside the Python runtime. An example is when I/O operations happen. I/O operations make lower level system calls in the OS that is outside of the Python runtime, and the GIL is released since we are not working directly with Python objects. The GIL is reacquired when the data is received and translated back into a Python object instance. In other languages such as Java, C, C++, the GIL is non-existent and threads running on different cores are parallel. *asyncio* exploits this fact that I/O releases GIL to give us concurrency with a single-thread using objects known as coroutines. A coroutine itself can be thought of as an even lighter-weight version of a thread. Of course, single-threads belong to the same process, so coroutines share the same memory space. To understand this single-threaded event loop model, we introduce sockets. A socket is a low-level abstraction for sending and receiving data over between applications over a network.

Sockets are by default blocking in nature, meaning that when we fire off a request for a data packet, it stops/blocks until we receive the data. However, they can be made to operate in a non-blocking mode, which means we fire and forget, until the OS notifies us that we have an inbox. These notification systems is interfaced by the OS and is OS-dependent. *asyncio* abstracts away these details to work with the different notification systems, allowing the developers to be mostly ignorant about the underlying mechanisms of data transfer. In *asyncio*, we we hit an I/O operation, we request the OS to make the I/O request, and then register in the notification system to remind us of any updates. Meanwhile, we are free to return to the Python runtime and execute other work. When the notification is received, Python receives the object and continues the code execution. The tasks that are waiting for I/O are tracked inside the Python runtime using a construct known as the event loop.

In *asyncio*, the event loop keeps a queue of tasks, which are wrappers around coroutines. When an event loop is instantiated, an empty queue of tasks is created. Each iteration of the event loop checks for tasks to run and runs them until an I/O operation is hit, which *asyncio* then makes the necessary system calls for the operation and for being notified of their progress. It then goes on to look for other tasks to execute. Each iteration also involves checking for any tasks that have already completed.

Figure 16.1: An example of a thread submitting tasks to the event loop.

Source: Figure 1.9 in Fowler [7]



Creating a coroutine is fairly straightforward, using the *async def* keyword and *await* keyword at the I/O code points. We will demonstrate this later in some substantial examples. When a coroutine is called in the manner *result = coroutine(args)*, it is not executed immediately, but a coroutine object is assigned to the variable *result*. To run a coroutine, it needs to be explicitly put on the event loop, and an iteration of the event loops need to be triggered. The different ways of doing so would be of primary interest going forward in this section.

One of the convenience functions is *asyncio.run*, we takes in a coroutine and runs it.

```

1 import asyncio
2
3 async def main():

```

```

4     await asyncio.sleep(1)
5     print("slept one(s)")
6     return True
7
8 if __name__ == "__main__":
9     result = asyncio.run(main())

```

Here, *asyncio* creates a new event-loop, and runs the *main* coroutine until completion. *asyncio.sleep* is also coroutine, which we will use as a stand-in function to simulate an I/O operation such as a database request. After completion, some cleanup is done and the event loop is closed. The *asyncio.run* function is intended to be the main entry into an asynchronous application. Other coroutines should be launched inside the coroutine that is called inside this *run* function, which is in our case, *main*. The *await* keyword is followed by an object of the *awaitable* type, which is satisfied by coroutines. The coroutine itself is paused on the *await sleep* instruction, and is then woken up later to continue execution. While no concurrent work was done here, these are the basic building blocks of an asynchronous application.

In fact, suppose we write this:

```

1 import asyncio
2
3 async def main():
4     await asyncio.sleep(1)
5     print("slept one(s)")
6     await asyncio.sleep(1)
7     print("slept one(s)")
8     return True
9
10 if __name__ == "__main__":
11     result = asyncio.run(main())

```

it will actually take roughly two seconds. As mentioned, the coroutine itself pauses until the line completes, which means we are sleeping one second twice, one after the other, instead of together. But we really want concurrency. In order to achieve this, we need to introduce *tasks*, which are wrappers around coroutines to schedule the contained coroutine to run on the event loop as early as possible. This scheduling and execution is non-blocking, while the *await* keyword is blocking. We can do this with *asyncio.create_task* function. When a task is created, it is scheduled to run on the next iteration of the event loop. An iteration of the event loop can be triggered by the *await* statement. Usually, when a task is created, we would await it at some point. An interesting aside is made. If we write this:

```

1 async def verbose_sleep(s):
2     await asyncio.sleep(s)
3     print("slept for", s)
4
5 async def main():
6     task1 = asyncio.create_task(verbose_sleep(5))
7     task2 = asyncio.create_task(verbose_sleep(6))
8     return True
9
10 if __name__ == "__main__":
11     asyncio.run(main())

```

the application would actually terminate immediately. Here the task objects are created but the event loop is not triggered and the function terminates. If instead we go with

```

1 async def verbose_sleep(s):
2     await asyncio.sleep(s)

```

```

3     print("slept for", s)
4
5 async def main():
6     task1 = asyncio.create_task(verbose_sleep(5))
7     task2 = asyncio.create_task(verbose_sleep(3))
8     await task2
9     return True
10
11 if __name__ == "__main__":
12     asyncio.run(main())

```

then we will see 'slept for 3' - both tasks are scheduled and run on the event loop, but the application terminates after the second task is completed. There are still more seconds to sleep on the first task and it never gets to finish. Instead,

```

1 async def verbose_sleep(s):
2     await asyncio.sleep(s)
3     print("slept for", s)
4
5 async def main():
6     task1 = asyncio.create_task(verbose_sleep(5))
7     task2 = asyncio.create_task(verbose_sleep(6))
8     await task2
9     return True
10
11 if __name__ == "__main__":
12     asyncio.run(main())

```

this would see both 'slept for 5' and 'slept for 6'. Even though we never awaited the first task, it runs and completes before the second task completes. Of course, if we await both tasks, we will always see both print statements. Tasks are scheduled to run as soon as possible, which generally means when the first *await* statement is encountered after the task has been created. The sleeping here is done concurrently, and in the last example, should terminate after approximately six seconds and not eleven. Of course, we may do some work without wasting time. Before we demonstrate that, let's create a utility decorator function that acts as a stopwatch timer for our application. We will discuss decorators in other sections

-

```

1 def asyn_timefn(func):
2     import time
3     from functools import wraps
4     @wraps(func)
5     async def timediff(*args, **kwargs):
6         a = time.time()
7         result = await func(*args, **kwargs)
8         print(f"@timefn: {func.__name__} took {time.time() - a} seconds")
9         return result
10    return timediff

```

Now suppose we want to compute the mean of ten million numbers of a uniform random sample.

```

1 async def verbose_sleep(s):
2     await asyncio.sleep(s)
3     print("slept for", s)
4
5 @asyn_timefn
6 async def main():

```

```

7     task1 = asyncio.create_task(verbose_sleep(5))
8     task2 = asyncio.create_task(verbose_sleep(6))
9     import random
10    import numpy as np
11    work = np.mean([random.uniform(0,1) for _ in range(10000000)])
12    await task1
13    await task2
14    return work
15
16 if __name__ == "__main__":
17     work = asyncio.run(main())
18     print(work)
19
20 '''
21 slept for 5
22 slept for 6
23 @timefn: main took 8.519409894943237 seconds
24 0.4998836811668402
25 '''

```

Interestingly it took quite a bit more than six seconds to complete this - it turns out that the mathematical computation was not done concurrently. This is because the tasks were created and scheduled but no event loop iteration was triggered. We can change the order of line 11 and 12, but this would only mean we are doing the computation and task2 in parallel. We want to do the math while waiting for both task1 and task2. Here's the trick:

```

1 @asyn_timefn
2 async def main():
3     task1 = asyncio.create_task(verbose_sleep(5))
4     task2 = asyncio.create_task(verbose_sleep(6))
5     import random
6     import numpy as np
7     await asyncio.sleep(0)
8     work = np.mean([random.uniform(0,1) for _ in range(10000000)])
9     await task1
10    await task2
11    return work
12 '''
13 slept for 5
14 slept for 6
15 @timefn: main took 6.001631021499634 seconds
16 0.49996590071870267
17 '''

```

The zero-sleep operation triggers the event loop, fires both sleep requests and does the mathematical work. There is almost no additional time spent on the statistical sampling work. Now we know how to work with tasks and coroutines, and some quirks. Let's discover some more of the *asyncio* workflow and control.

Network connections can be unstable and might hang indefinitely. We want to be able to cancel tasks. *Task* objects come with method *cancel*, and the running task raises a *CancelledError* when we await on the task.

```

1 async def main():
2     task1 = asyncio.create_task(verbose_sleep(5))
3     assert not task1.done()

```

```

4     task1.cancel()
5     try:
6         await task1
7     except asyncio.CancelledError:
8         print("task was cancelled")
9     return
10
11 if __name__ == "__main__":
12     asyncio.run(main())
13
14 '''
15 task was cancelled
16 '''

```

It should be noted that the *CancelledError* is only thrown at *await* statements, meaning to say if the cancel is submitted while the task is midway through executing plain Python code, the code runs until the next *await* statement, if any. We might already have some time in mind to cancel the task, which we should then go for *asyncio.wait_for*, taking in a coroutine/task and number of seconds to wait for - after which a *TimeoutException* is thrown:

```

1 async def main():
2     try:
3         task1=asyncio.create_task(verbose_sleep(3))
4         await asyncio.wait_for(task1,2)
5         await task1
6     except asyncio.TimeoutError:
7         print("task was timed out with status cancel =", task1.cancelled())
8     return
9
10 if __name__ == "__main__":
11     asyncio.run(main())
12
13 '''
14 task was timed out with status cancel = True
15 '''

```

But suppose we want to let the task continue executing, instead of cancelling it after specified seconds we can *asyncio.shield* the function:

```

1 async def main():
2     try:
3         task1=asyncio.create_task(verbose_sleep(3))
4         await asyncio.wait_for(asyncio.shield(task1),2)
5         await task1
6     except asyncio.TimeoutError:
7         print("task was timed out with status cancel =", task1.cancelled())
8         await task1
9     return
10
11 if __name__ == "__main__":
12     asyncio.run(main())
13
14 '''
15 task was timed out with status cancel = False
16 slept for 3
17 '''

```

Ok now let's understand what links tasks and coroutines. We know they are both awaitable, but what is an awaitable? To do that we need to introduce *futures*. Futures are Python objects that contain a single value you expect to get in the future. When you first create one, that value normally doesn't exist, and somewhere downstream, you have it, and you set the value of the future to that object. Futures can be awaited, and are either *done* or *undone*. The following example should be sufficiently clear:

```
1 async def main():
2     future = asyncio.Future()
3     print(future.done())
4     async def _helper():
5         await asyncio.sleep(3)
6         future.set_result("DONE")
7
8     task = asyncio.create_task(_helper())
9     await asyncio.sleep(0)
10    value = await future
11    print(future.done(), "with result", value)
12    return
13
14 if __name__ == "__main__":
15     asyncio.run(main())
16 '''
17 False
18 True with result DONE
19 '''
```

The internal implementation of the *asyncio* API relies heavily on futures. A task can be thought of as combination of coroutine and future. When a task is created, an empty future is created - with the expectation that the future result is set to the output of the coroutine in question.

They are awaitables, which means to say they are object instances of classes that implement the *Awaitable* abstract base class. Abstract base classes are discussed in other sections, but briefly, it means they specify some behavior that parent classes implementing it should satisfy. It determines the contract specification. For instance, when we hire a burger cook, when we call *chef.cook()*, we know we want a burger, even though different cooks can prefer sous vide, grill or whatever. In our case the method to be implemented is the *__await__* method. The inheritance diagram looks something like this:

```
Awaitable
    Coroutine
    Future
        Task
```

When we look at the API documentation for *asyncio*, it is useful to know what is the object instance they require. The children is always an instance of the parent class, so for a function that takes in an awaitable, we really can pass in an instance of coroutine, future or task. A function that takes in instances of coroutine only would not work when a task is passed in, however.

Common mistakes in using the *asyncio* API are discussed. One is trying to use CPU-bound code in tasks or coroutines without multiprocessing. See that the following two listings that approximately the same time.

```
1 async def cpu_work():
2     import numpy as np
3     import random
```



```

4     sample = [random.uniform(0,1) for _ in range(10000000)]
5     sum = 0
6     for num in sample:
7         sum += num
8     return sum
9
10 @asyn_timefn
11 async def main():
12     task1 = asyncio.create_task(cpu_work())
13     task2 = asyncio.create_task(cpu_work())
14
15     await task1
16     await task2
17     return
18
19 if __name__ == "__main__":
20     asyncio.run(main())
21 '''
22 @timefn: main took 4.500612020492554 seconds
23 '''

```

```

1 def cpu_work():
2     import numpy as np
3     import random
4     sample = [random.uniform(0,1) for _ in range(10000000)]
5     sum = 0
6     for num in sample:
7         sum += num
8     return sum
9
10 @asyn_timefn
11 async def main():
12     cpu_work()
13     cpu_work()
14     return
15 '''
16 @timefn: main took 4.320184946060181 seconds
17 '''

```

Actually the second program runs faster, because there is lesser overhead without coroutines, Although the difference is marginal, it definitely does not help the program. Firstly, the CPU work here does not have any await statements. So it would not give up execution for other I/O tasks behind it in the event loop to get a chance to submit the I/O request until it finishes execution. Secondly, because *asyncio* is single threaded, no CPU work occurs in parallel. Furthermore, even if it were multi-threaded, the GIL would prevent any parallelism. What we would want here is multi-processing, which would be discussed later. The second mistake would be to asynchronously call blocking APIs. When we know that we are making an I/O operation through an API, we might be tempted to wrap it in a coroutine. But the APIs themselves may be blocking the event loop itself, so we will not get performance benefits.

```

1 async def hangukquant():
2     import requests
3     res=requests.get("https://www.hangukquant.com")
4     return
5
6 @asyn_timefn

```

```

7 async def main():
8     tasks=[asyncio.create_task(hangukquant()) for _ in range(30)]
9     await asyncio.sleep(0)
10    [await task for task in tasks]
11    return
12
13 if __name__ == "__main__":
14     asyncio.run(main())
15 '''
16 @timefn: main took 6.972156286239624 seconds
17 '''

```

That the network request is an I/O operation is accurate. However, the internal implementation of the *requests* module block on the *get* method, so no asynchronous work is done here. We can either use multithreading with blocking APIs such as *requests* or use dedicated asyncio-supported libraries such as *aiohttp*, which uses non-blocking sockets. Later, we will examine them.

Previously we used the *asyncio.run* method to handle the event loop for us. We can actually have more fine grained control over the event loop behavior, but we would need to be careful about cleaning up the resource and closing the loop. This would work:

```

1 async def main():
2     await verbose_sleep(1)
3     return
4
5 if __name__ == "__main__":
6     loop = asyncio.new_event_loop()
7     try:
8         loop.run_until_complete(main())
9     finally:
10        loop.close()

```

We can also replace *new_event_loop* with *get_event_loop*, which gets the running loop if any, or creates a new one. If we do *get_running_loop* instead, we would get an error that there is no currently running event loop. We can also explicitly put non-async functions on the event loop to run on the next iteration with *call_soon*:

```

1 def print_hello():
2     print("hello")
3
4 async def main():
5     loop = asyncio.get_running_loop()
6     loop.call_soon(print_hello)
7     await verbose_sleep(1)
8     return
9
10 if __name__ == "__main__":
11     loop = asyncio.new_event_loop()
12     try:
13         loop.run_until_complete(main())
14     finally:
15         loop.close()
16 '''
17 hello
18 slept for 1
19 '''

```

We refer readers to documentation to nuances, if more control over the event loop is desired within the application. Other options, such as debug mode exists, i.e. `asyncio.run(main(), debug=True)`. When working with production quality applications, exception handling is very important. Note that when exception is thrown inside a running task, the task is considered done with the result as exception. No exception is thrown up the call stack, and there is no cleanup. To get the exception thrown to us, we need to *await* the task, otherwise the exception may never be retrieved, particularly if the task is not garbage collected.

We want to have control over the workflow of coroutines, and explore some useful APIs that are *asyncio* compatible. An example is the *aiohttp* library, which uses non blocking sockets to make web requests and returns coroutines for the requests. As mentioned, the *requests* library is blocking and would block the thread it runs in. Before we go into *aiohttp*, let us first understand the role of context managers in resource control. In particular, we are interested in asynchronous context managers - they implement two coroutines `__aenter__` and `__aexit__`. Among other things, context managers help to manage resources and perform auxiliary functionalities such as resource clean up. A common example would be:

```
1 with open("test.txt") as file:
2     lines= file.readlines()
```

which is a more Python-ic way of doing this:

```
1 file = open("test.txt")
2 try:
3     lines= file.readlines()
4 finally:
5     file.close()
```

With asynchronous context managers, the syntax is slightly different. For instance, we will use *async-with* instead of the *with* keyword. We can acquire an asynchronous lock like this:

```
1 lock = asyncio.Lock()
2 counter = 0
3 async def increment():
4     global counter
5     async with lock:
6         counter += 1
```

We will not go into too intimate details with context managers - they will be discussed in other sections. *aiohttp* uses async context managers for acquiring HTTP sessions and connections. Sessioning is a networking term for keeping connections open. Without going into the details, when HTTP requests are made over the network, protocols such as the TCP handshake and exchange of other information are expensive. Sessioning keeps the connections open and allow us to recycle them, and this is known as connection pooling, which reduces the resources required and improve the performance of *aiohttp*-based applications. With this non-blocking library, we may use sessioning and coroutines to speed up the same problem in Listing 16.1.1.

```
1 async def hangukquant(session):
2     async with session.get("https://www.hangukquant.com") as res:
3         return res.text
4
5 @asyn_timefn
6 async def main():
7     import aiohttp
8     async with aiohttp.ClientSession() as session:
9         tasks=[asyncio.create_task(hangukquant(session)) for _ in range(30)]
```

```

10     await asyncio.sleep(0)
11     [await task for task in tasks]
12     return
13
14 if __name__ == "__main__":
15     asyncio.run(main())
16 '''
17 @timefn: main took 1.417881965637207 seconds
18 '''

```

Notice that our program runs significantly faster. The *ClientSession* creates a default maximum of hundred connections. This may be altered - we refer readers to *aiohttp* documentation. Previously, we could set timeouts with *asyncio.wait_for* method. Although this works with arbitrary coroutines, and hence also applicable here, *aiohttp* provides out of the box functionalities to do this cleanly. At the session level, we can instantiate the client session with a timeout object to get a timeout limit and replace the relevant lines of code with:

```

1     session_timeout = aiohttp.ClientTimeout(total=1,connect=0.3)
2     timed_session = aiohttp.ClientSession(timeout=session_timeout)
3     async with timed_session:
4         ...

```

which sets a time limit on the whole request and also for subroutines such as establishing connections. On the per-request level, we can also specify timeouts by passing a client timeout object to the get request:

```

1     async with session.get(
2         "https://www.hangukquant.com",
3         timeout=aiohttp.ClientTimeout(total=0.5)
4     ) as res:
5         ...

```

If the requests are not successful by the time specified, we would get an *asyncio.TimeoutError* on await statement. The last code example works, but it is kind of clunky. First, we have to wrap coroutines in task objects, then trigger an iteration of the event loop with a pseudo-sleep method and then await each task. We can do this cleanly using the *asyncio.gather* method, which takes in a sequence of awaitables and runs them concurrently. If any of the awaitables passed in were coroutines, they are automatically wrapped in task objects so that they get run on the event loop concurrently. The gather method itself returns an awaitable, and even though the individual tasks may not complete in deterministic order, the results are returned in the order we pass them in. Refer to the following example:

```

1 async def hangukquant(session, request_no):
2     async with session.get("https://www.hangukquant.com") as res:
3         result=res.text
4     return request_no
5
6 async def main():
7     import aiohttp
8     async with aiohttp.ClientSession() as session:
9         results = await asyncio.gather(*[hangukquant(session,i) for i in range(30)])
10        print(results)
11    return
12
13 if __name__ == "__main__":
14     asyncio.run(main())

```

```

15 '''
16 [0, 1, 2, 3, ..., 28, 29]
17 '''

```

Much cleaner. What about exceptions? *gather* has an optional parameter *return_exceptions*, specifying two different behavior - it defaults to *False*, which causes *gather* to throw an exception when we await it if any of the coroutines throw an exception. It should be noted that in these scenarios, the other coroutines are not actually cancelled - if we handle the exception, they will be allowed to run until completion. Otherwise, they will run until they are cleaned up downstream or to completion, whichever is earlier. If we set *return_exceptions* to *True*, then the *gather* method will not throw any exception. Instead, the return value will be thrown exception inside the coroutine:

```

1 async def hangukquant(session, request_no):
2     async with session.get("https://www.hangukquant.com") as res:
3         import random
4         if random.uniform(0,1) > 0.5:
5             raise Exception("unlucky error")
6         result=res.text
7         return request_no
8
9 @async_timefn
10 async def main():
11     import aiohttp
12     async with aiohttp.ClientSession() as session:
13         results = await asyncio.gather(*[hangukquant(session,i) for i in range(10)],
14         return_exceptions=True)
15         print(results)
16         return
17
18 if __name__ == "__main__":
19     asyncio.run(main())
20 '''
21 [0, Exception('unlucky error'), 2, Exception('unlucky error'), Exception('unlucky error'),
22 Exception('unlucky error'), Exception('unlucky error'), Exception('unlucky error'),
23 8, 9]
24 '''

```

Note that if let the exception be thrown at *gather*, then we will only retrieve the first one when we await it. Maybe this is okay, but maybe we do want to specify a different behavior. Perhaps we would like to cancel our tasks if one of the tasks fail, such as when the tasks make API requests to the same server and we receive a 429 code as rate limits get triggered. Spamming the server with the remaining requests might get us blacklisted - so cancelling the remaining tasks would be desirable. If the tasks take disproportionate amount of time to complete, the *gather* method only completes after the longest task is done - we may prefer to work on the results of the requests that have already completed first. *asyncio* provides the *as_completed* method which takes in a list of awaitables and returns an iterator of futures. We can iterate over and await on each of the futures, with the result being the return value of the coroutine on a first-finish basis. Therefore, the ordering of results are non-deterministic. We demonstrate with the following:

```

1 async def work(num):
2     import random
3     work_cost = random.uniform(0,4)

```

```

4     await asyncio.sleep(work_cost)
5     return (num,work_cost)
6
7 @asyn_timefn
8 async def main():
9     async_work = [work(i) for i in range(10)]
10    for finished in asyncio.as_completed(async_work):
11        res = await finished
12        print(res)
13    return
14
15 if __name__ == "__main__":
16     asyncio.run(main())
17
18 '''
19 (6, 0.3673681494715675)
20 (9, 1.276757248018661)
21 (7, 1.3913958197279603)
22 (3, 2.3091369741247556)
23 (5, 2.3978259114992753)
24 (8, 2.4415577155094477)
25 (2, 2.5467792505773)
26 (1, 2.60398276135025)
27 (0, 3.159090865296541)
28 (4, 3.824869021632983)
29 @timefn: main took 3.829425096511841 seconds
30 '''

```

Since we have access to individual tasks, we also have better control over exception handling. As before, any task exception would be thrown at the `await` statement. There is an optional parameter `timeout` in the `asyncio.as_completed` function, which specifies the number of seconds to let the group of tasks run. Any task taking longer than the timeout would throw a `asyncio.TimeoutException` when awaited. As we get the results of those that finished earlier right away, any result retrieved may straight away be worked on without the results of the others. While this is an advantage over `gather`, we might want even more control. When timeout occurs, exception is thrown, but there is no cancellation. The tasks still run in the background. Perhaps, we would like to cancel them. `asyncio` offers another method `wait` that gives several options on when we would like to receive our results from the tasks. The method returns two sets - one in which the tasks are completed (either a result is computed or exception is thrown) and the other in which the tasks are still running. The optional `timeout` configuration on this method specifies the time requirements. Additionally, upon timeout, no exception is thrown, the coroutines are not cancelled and we would have to iterate over the pending task set to cancel them. Instead of throwing the exception, on timeout, all the tasks still running are just returned in the pending set for us to do our preferred handling. The options in the `wait` protocol is `ALL_COMPLETED`, `FIRST_EXCEPTION`, `FIRST_COMPLETED`. It should be pretty clear from their variable names what behavior they specify. We refer the details to the documentation, and instead provide an example. To be brief, the `ALL_COMPLETED` option would behave similar to `gather`, and the `FIRST_COMPLETED` option would behave similar to `as_completed`. Of course, if we receive the results of the tasks that completed earlier, we may go ahead and perform work on it without waiting for the other tasks to complete. `FIRST_EXCEPTION` option is interesting - we shall demonstrate a program that cancels the pending tasks on encountering an exception:

```

1 async def work(num):

```

```

2     import random
3     work_cost = random.uniform(0,1)
4     await asyncio.sleep(work_cost)
5     if random.uniform(0,1) < 0.30:
6         raise Exception(f"Task {num} had an unlucky failure")
7     return (num,work_cost)
8
9 @asyn_timefn
10 async def main():
11     async_work = [asyncio.create_task(work(i)) for i in range(10)]
12     done, pending = await asyncio.wait(async_work, return_when=asyncio.FIRST_EXCEPTION)
13     for item in done:
14         if item.exception() is None:
15             print(item.result())
16         else:
17             print(item.exception())
18     for item in pending:
19         print("cancelling ", async_work.index(item))
20         item.cancel()
21     return
22
23 if __name__ == "__main__":
24     asyncio.run(main())
25
26 '''
27 Task 7 had an unlucky failure
28 (4, 0.3565076327157949)
29 (2, 0.20337667536574278)
30 (9, 0.1277640138373286)
31 cancelling 5
32 cancelling 8
33 cancelling 3
34 cancelling 6
35 cancelling 0
36 cancelling 1
37 '''

```

Note that the *done* variable set is guaranteed to contain minimally one element - the task that threw an exception is considered as done.