

Bài 52 - Parallel Computing on Python

30 Nov 2020 - phamdinhhkhanh

Menu

- 1. Tại sao cần xử lý song song
 - 1.1. Thread và Process
- 2. Khởi tạo thread trong python
 - 2.1. Khởi tạo từ hàm
 - 2.2. Khởi tạo kế thừa
 - 2.3. Cơ chế Thread Lock
- 3. Khởi tạo process trong python
 - 3.1. Sử dụng chung dữ liệu
 - 3.2. Shared object giữa các Process
 - 3.3. Pool trong multiprocessing
- 4. Process Pool vs Thread Pool
 - 4.1. Process Pool
 - 4.2. Thread Pool
- 5. Queue
- 6. Kết luận
- 7. Tài liệu

1. Tại sao cần xử lý song song

Trong quá trình xây dựng các ứng dụng deep learning trên python, mình nhận ra rằng để tạo ra một ứng dụng thì không khó. Nhưng để tạo ra một ứng dụng đáp ứng được tốc độ xử lý, độ chính xác và mức độ sử dụng resource thì cần phải tối ưu rất nhiều thứ. Bạn sẽ phải quan tâm đến các khía cạnh như:

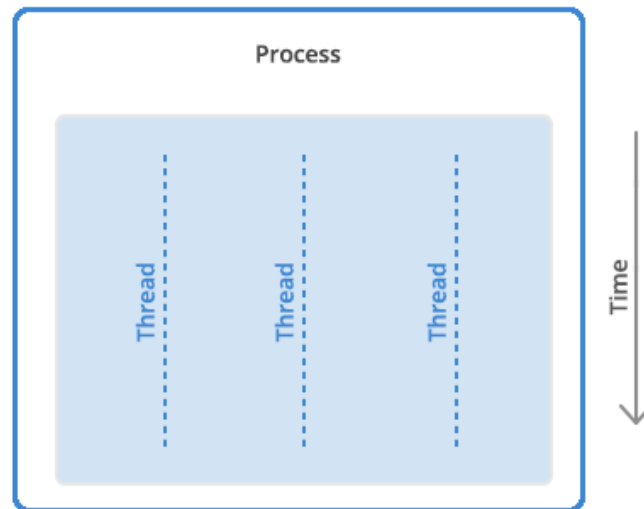
- Giảm nhẹ kích thước mô hình thông qua: Quantization và compress mô hình.
- Optimize lại code.
- Chuyển từ single-thread sang multi-thread.
- Allocate lại tài nguyên như CPU, Memory.

Đặc biệt là các ứng dụng trên python thì tối ưu tốc độ xử lý là một challenge bởi python bị ràng buộc bởi cơ chế GIL (Global Interpreter Lock). Tức là nó chỉ cho phép một thread hoạt động truy suất và chỉnh sửa bộ nhớ tại một thời điểm. Do đó python không tận dụng được các tính toán đa luồng. Tuy nhiên ở python 3.2 trở đi thì python đã bắt đầu hỗ trợ đa luồng. Và thông qua bài viết này mình sẽ hướng dẫn các bạn có thể accelerate các ứng dụng của mình thông qua đa luồng.

Nhưng trước tiên chúng ta cần hiểu về thread/process là gì? Vì blog dành cho đa dạng bạn đọc ở trình độ khác nhau nên bạn nào đã biết thì có thể bỏ qua phần kiến thức rất sơ đẳng này.

1.1. Thread và Process

Top



Thread và process là hai khái niệm cơ bản trong lập trình và cũng có nhiều định nghĩa từ các nguồn khác nhau cho chúng.

Process là gì?

Chúng ta hiểu một cách đơn giản thì process là tiến trình để chạy một phần mềm. Khi bạn start một program thì tức là bạn đang khởi tạo một process. Hệ điều hành khi đó sẽ cung cấp các tài nguyên về memory, cpu, disk, bandwidth cho process để cho chạy ứng dụng của bạn.

Task Manager					
File Options View					
Processes Performance App history Startup Users Details Services					
Name	Status	6% CPU	43% Memory	15% Disk	0% Network
Apps (3)					
> Google Chrome (32 bit)		0.1%	107.8 MB	0.1 MB/s	0 Mbps
> Snagit Editor (32 bit)		0%	27.2 MB	0 MB/s	0 Mbps
> Task Manager		0.4%	12.9 MB	0.1 MB/s	0 Mbps
Background processes (58)					
> Andrea filters APO access servic...		0%	0.9 MB	0 MB/s	0 Mbps
AutoHotkey Unicode 64-bit		0.1%	2.3 MB	0 MB/s	0 Mbps
CCleaner		0%	16.5 MB	0 MB/s	0 Mbps
CCleaner		0%	4.6 MB	0 MB/s	0 Mbps
COMODO Internet Security		0%	0.6 MB	0 MB/s	0 Mbps
COMODO Internet Security		0%	0.8 MB	0 MB/s	0 Mbps
COMODO Internet Security		0%	0.4 MB	0 MB/s	0 Mbps
> COMODO Internet Security		0%	18.8 MB	0.1 MB/s	0 Mbps
Cortana		0%	0.1 MB	0 MB/s	0 Mbps
<div> ^ Fewer details Top </div>					

Hình 1: Khi bạn vào task management của window bạn có thể theo dõi các process đang chạy với mã PID của process. Mỗi một process sẽ phụ trách một instance của OS system và được cung cấp các thành phần như memory, cpu, disk, bandwidth,....

Lịch xử lý của các processes sẽ được OS sắp xếp dựa trên một số thuật toán lập lịch như round robin, first come first serve,... Mình sẽ không đi sâu vào phần này, các bạn có thể tham khảo thêm Operating System Scheduling algorithms (https://www.tutorialspoint.com/operating_system/os_process_scheduling_algorithms.htm).

Threads là gì?

Chắc hẳn bạn đã từng nghe đến thông số số cores của CPU. Các CPU càng hiện đại, số lượng cores sẽ càng nhiều. Các core sẽ hỗ trợ cho việc tính toán multi-task tốt hơn. Các threads sẽ được vận hành và tính toán trên các core của CPU.

Một process khi được khởi tạo sẽ sinh ra các threads để run application. Bạn sẽ thắc mắc vậy thì chỉ cần một thread cũng được ? Tại sao lại cần nhiều threads? Nhiều threads sẽ giúp cho việc tính toán multi-task tốt hơn. Tức là bạn có thể làm nhiều nhiệm vụ một lúc. Nếu coi mỗi thread là một công nhân, thì việc sản xuất sẽ nhanh hơn nếu có nhiều công nhân phối hợp cùng làm việc. Bạn có thể hình dung dễ hơn qua ví dụ:

Khi bạn làm việc với microsoft word, bạn gõ bàn phím thì có những công việc sau cần thực hiện:

- Đọc thông tin input từ keyboards.
- Hiển thị lên màn hình các thông tin đã nhập trong quá trình gõ.
- Highlight những chỗ bị sai chính tả.
- Suggest các từ để có thể gõ nhanh hơn.

Mỗi công việc được phụ trách bởi một thread và chúng phối hợp với nhau để giúp ứng dụng của bạn mượt hơn. Nếu chỉ có một thread làm tất cả mọi công việc thì nó sẽ bị quá tải và bạn có thể gặp phải giới hạn về tốc độ xử lý của CPU, thuật ngữ hay được gọi là CPU bound.

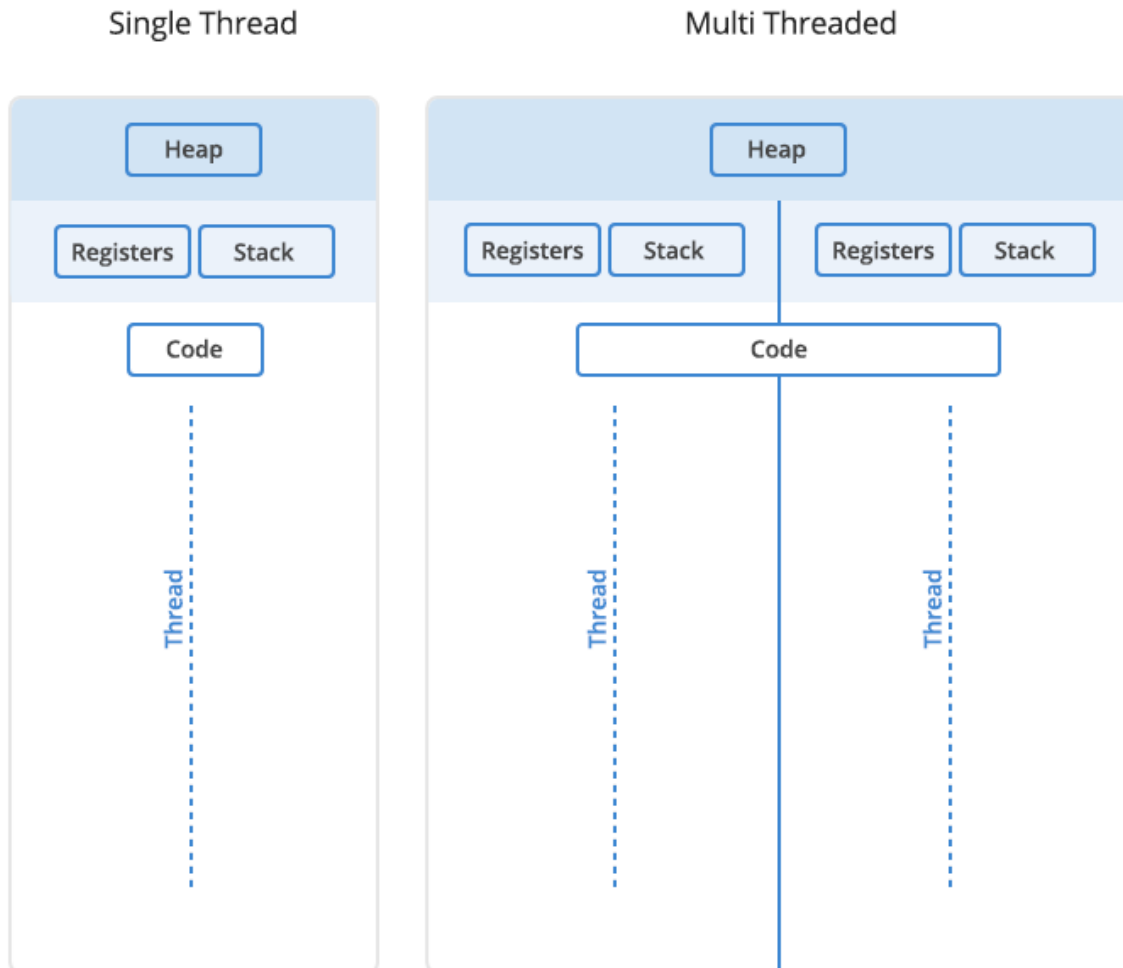
Một process có thể là single-thread hoặc multiple-threads tùy thuộc vào số lượng là một hoặc nhiều. Khi có nhiều threads thì đòi hỏi phải có sự phối hợp tính toán song song (parallel computing) giữa các threads với nhau. Từ đó sinh ra các khái niệm về đồng bộ (synchronous) và bất đồng bộ (asynchronous). Chúng ta sẽ làm rõ hai khái niệm này ở các phần tiếp theo.

Bạn có thể thắc mắc multiple-threads thì có khác gì khác biệt so với việc sử dụng multiple-processes? Chúng ta vẫn có thể tính toán song song được trên cả hai? Vậy tại sao lại cần phải tách một process thành nhiều threads làm gì ? Thực tế là trong python thì process và thread cùng kế thừa chung một interface là một base thread. Chúng sẽ có những đặc tính chung, nhưng thread là một phiên bản nhẹ hơn so với process. Do đó việc khởi tạo thread sẽ nhanh hơn. Một điểm khác biệt nữa đó là thread được thiết kế để có thể hoạt động tương tác lẫn nhau. Các threads trong cùng một process sẽ chia sẻ được dữ liệu qua lại nên có lợi thế về I/O. Dữ liệu của process thì được thiết kế private nên một process không thể chia sẻ dữ liệu với các process khác. Đây là lý do chúng ta cần nhiều threads hoạt động trong một process.

Tiếp theo môi trường hoạt động của multiple-threads sẽ như thế nào ?

Khi các threads chạy song song trên cùng một process, chúng sẽ khởi tạo dữ liệu như thế nào? Dữ liệu sẽ được lưu vào đâu? Chúng chia sẻ chung một code như thế nào? Chúng ta cùng làm rõ qua sơ đồ bên dưới.

Top



Hình 2: Cấu trúc của single thread và multiple threads.

Đầu tiên ứng dụng của bạn sẽ khi khởi chạy sẽ load code lên. Phần main của chương trình sẽ được compiler khởi chạy đầu tiên. Lần lượt các method sẽ được load vào môi trường stack theo trình tự chạy. Compiler chạy lần lượt các hàm trong stack. Các hàm được compiler biên dịch thành mã máy (byte code) và được thực thi để sinh ra dữ liệu. Dữ liệu sau đó được lưu trữ tại hai bộ nhớ là Heap và Stack (cái này cũng tùy thuộc vào virtual machine của từng ngôn ngữ). Stack lưu trữ method và các local variable còn heap lưu trữ object, array từ chương trình của bạn (phần lưu trữ này cũng có thể thay đổi tùy vào cách sắp xếp bộ nhớ của các ngôn ngữ). Nếu load trên stack thì không bị phân mảnh dữ liệu và có thời gian load/access nhanh hơn. Còn heap sẽ allocate vùng nhớ ngẫu nhiên, các ô nhớ không liên tục nên do đó bị phân mảnh.

Okie, mình nghĩ lý thuyết như vậy là đủ rồi. Tiếp theo chúng ta sẽ cùng thực hành khởi tạo các thread và process trong python.

2. Khởi tạo thread trong python

2.1. Khởi tạo từ hàm

Trên python3 để khởi tạo một thread thì chúng ta sử dụng module `_thread` , trên python2 là `thread` . Để start một method trên thread thì chúng ta chỉ cần truyền vào `_thread.start_new_thread()` tên method và các giá trị đối số của nó. Ví dụ bên dưới chúng ta sử dụng hàm `_counter()` để đếm lùi các số từ trên xuống dưới.

Top

```

1  import _thread
2  import time
3
4  def _counter(counter, thread_name):
5      while (counter):
6          time.sleep(0.01)
7          print("{}: {}".format(thread_name, counter))
8          counter -= 1
9
10 counter = 5
11
12 # Khởi tạo 2 threads 1 và 2
13 try:
14     _thread.start_new_thread(_counter, (counter, "khanh thread")) # pas
15     print("\n")
16     _thread.start_new_thread(_counter, (counter, "ai thread"))
17 except:
18     print("Error: unable to start thread")
19
20 # Running counter
21 while (counter):
22     counter -= 1
23     pass

```

```

1  ai thread: 5
2  khanh thread: 5
3  ai thread: 4
4  khanh thread: 4
5  ai thread: 3
6  khanh thread: 3
7  ai thread: 2
8  khanh thread: 2
9  ai thread: 1
10 khanh thread: 1

```

Ta thấy hai thread đã xen kẽ nhau cùng thực hiện tác vụ đếm ngược. Tuy nhiên về bản chất thì chúng vẫn là đơn luồng vì cơ chế GIL của python ép buộc một thời điểm chỉ có một thread được tương tác với dữ liệu. Có khá nhiều developer tỏ ra thất vọng về hạn chế này nhưng một số khác thì bảo vệ quan điểm này bởi nó giúp một dữ liệu không bị sử dụng và thay đổi cùng lúc bởi nhiều threads. Hiện tượng này dẫn tới concurrency, một trong những bug thường gặp ở các ngôn ngữ đa luồng như java hay C++.

Trong ví dụ trên thì các method trên hai threads khanh và ai khởi chạy độc lập nhau mà không ưu tiên một thread hoàn thành thì mới chạy thread tiếp theo. Cách chạy như vậy được gọi là bất đồng bộ asynchronous, một khái niệm cơ bản của parallel application. Trái ngược lại thì là đồng bộ synchronous, các method sẽ chạy theo tuần tự, sau khi method trước đó đã hoàn thành.

2.2. Khởi tạo kế thừa

Một cách khác để khởi tạo một thread đó là kế thừa lại Threading module. Kiểu kế thừa này khá phổ biến trong lập trình, chắc các bạn còn nhớ khi khởi tạo model trên pytorch chúng ta cũng kế thừa lại nn.Module (<https://phamdinhhkhanh.github.io/2019/08/10/PytorchTutorial1.html#31->

ki%E1%BA%BFn-tr%C3%BAc-m%E1%BA%A1ng-cnn) chứ ? Khi đó chúng ta chỉ cần override lại các method cần điều chỉnh từ class cha.

```

1      import threading
2      import time
3
4      class FirstThread(threading.Thread):
5          def __init__(self, thread_id, thread_name, counter):
6              threading.Thread.__init__(self)
7              self.thread_id = thread_id
8              self.thread_name = thread_name
9              self.counter = counter
10
11         def run(self):
12             print("Start thread {}".format(self.thread_name))
13             while (self.counter):
14                 time.sleep(0.01)
15                 print("{} : {}".format(self.thread_name, self.counter))
16                 self.counter -= 1
17             print("End thread {}".format(self.thread_name))
18
19
20     thread1 = FirstThread(1, "khanh thread", 5)
21     thread2 = FirstThread(2, "ai thread", 5)
22
23     thread1.start()
24     thread2.start()

```

```

1      Start thread khanh thread!
2      Start thread ai thread!
3      khanh thread : 5
4      ai thread : 5
5      khanh thread : 4
6      ai thread : 4
7      khanh thread : 3
8      ai thread : 3
9      khanh thread : 2
10     ai thread : 2
11     khanh thread : 1
12     End thread khanh thread
13     ai thread : 1
14     End thread ai thread

```

2.3. Cơ chế Thread Lock

Như đã giới thiệu chương trước, trong ví dụ ở trên các threads là bất đồng bộ (asynchronous). Hai threads chạy độc lập với nhau mà không theo thứ tự. Chúng ta có thể đồng bộ (synchronous) các thread. Tức là cho phép một thread chạy xong thì thread khác mới được phép chạy bằng cách sử dụng Thread Lock trong python.

Top

```

1  import threading
2
3  class FirstThread(threading.Thread):
4      def __init__(self, thread_id, thread_name, counter):
5          threading.Thread.__init__(self)
6          self.thread_id = thread_id
7          self.thread_name = thread_name
8          self.counter = counter
9
10     def run(self):
11         threadLock.acquire()
12         print("Start thread {}".format(self.thread_name))
13         while (self.counter):
14             time.sleep(0.01)
15             print("{} : {}".format(self.thread_name, self.counter))
16             self.counter -= 1
17         print("End thread {}".format(self.thread_name))
18         threadLock.release()
19
20     threadLock = threading.Lock()
21     thread1 = FirstThread(1, "khanh thread", 5)
22     thread2 = FirstThread(2, "linh thread", 5)
23
24     thread1.start()
25     thread2.start()
26
27     threads = [thread1, thread2]
28
29     for t in threads:
30         t.join()

```

```

1  Start thread khanh thread!
2  khanh thread : 5
3  khanh thread : 4
4  khanh thread : 3
5  khanh thread : 2
6  khanh thread : 1
7  End thread khanh thread
8  Start thread linh thread!
9  linh thread : 5
10  linh thread : 4
11  linh thread : 3
12  linh thread : 2
13  linh thread : 1
14  End thread linh thread

```

Trong hàm `run()` của thread thì chỉ cần thêm hàm `thread.acquire()` và `thread.release()` vào đầu và cuối hàm thì luồng sẽ được locking cho đến khi thread chạy xong thì thread khác mới được xử lý tiếp. Như chúng ta thấy, sau khi thread1 xử lý xong thì mới đến lượt thread2 xử lý.

3. Khởi tạo process trong python

Để khởi tạo một process trong python chúng ta sử dụng class `Process` của thư viện `multiprocessing`. Chúng ta cũng truyền vào hàm và đối số như đã thực hiện với thread.

Top

```

1  from multiprocessing import Process
2  import time
3
4  def _counter(counter, process_name):
5      while (counter):
6          time.sleep(0.01)
7          print("{}: {}".format(process_name, counter))
8          counter -= 1
9
10 counter = 5
11
12 exec1 = Process(target=_counter, args=(counter, "khanh thread")) # pa.
13 exec2 = Process(target=_counter, args=(counter, "ai thread"))
14
15 exec1.start()
16 exec2.start()
17
18 for exec in execs:
19     exec.join()

```

```

1  khanh thread: 5
2  ai thread: 5
3  khanh thread: 4
4  ai thread: 4
5  khanh thread: 3
6  ai thread: 3
7  khanh thread: 2
8  ai thread: 2
9  khanh thread: 1
10 ai thread: 1

```

Khi làm việc với multi-process, chúng ta luôn cần một lệnh `join()` để đảm bảo main process hoàn thành sau cùng sau khi các child process khác kết thúc.

Ta nhận thấy rằng các process được thực hiện một cách độc lập và bất đồng bộ. Để đồng bộ các process với nhau thì chúng ta đơn giản là `lock` chúng lại.

Top


```

1  from multiprocessing import Process, Lock
2  import time
3
4  def _counter_lock(counter, process_name, lock):
5      lock.acquire()
6      while (counter):
7          time.sleep(0.01)
8          print("{}: {}".format(process_name, counter))
9          counter -= 1
10         lock.release()
11
12     counter = 5
13
14     lock = Lock()
15     exec1 = Process(target=_counter_lock, args=(counter, "khanh thread", lock))
16     exec2 = Process(target=_counter_lock, args=(counter, "ai thread", lock))
17     execs = [exec1, exec2]
18
19     for exec in execs:
20         exec.start()

```

```

1  khanh thread: 5
2  khanh thread: 4
3  khanh thread: 3
4  khanh thread: 2
5  khanh thread: 1
6  ai thread: 5
7  ai thread: 4
8  ai thread: 3
9  ai thread: 2
10 ai thread: 1

```

Bạn thấy đó, chúng cũng khá nán thread phải không nào ?

3.1. Sử dụng chung dữ liệu

Khi làm việc với các ứng dụng concurrent thì chúng ta nên hạn chế nhất có thể việc chia sẻ dữ liệu giữa các process để tránh xảy ra các lỗi phát sinh do concurrency. Tuy nhiên python vẫn cung cấp một cơ chế giúp chia sẻ dữ liệu giữa các process, đó chính là các shared memory object trong multiprocessing như Value, Array. Thật vậy, giả sử ở ví dụ bên dưới chúng ta sử dụng 2 processes để thay đổi dấu các phần tử của một list các số nguyên.

Top

```

1  from multiprocessing import Process, Lock
2  import time
3
4  def _counter_arr(arrs, process_name):
5      lock.acquire()
6      for i, el in enumerate(arrs):
7          time.sleep(0.01)
8          arrs[i] = -arrs[i]
9          print("{}: {}".format(process_name, arrs[i]))
10     lock.release()
11
12     arrs = [1, 2, 3, 4]
13     lock = Lock()
14     exec1 = Process(target=_counter_arr, args=(arrs, "khanh process")) #
15     exec2 = Process(target=_counter_arr, args=(arrs, "ai process"))
16     execs = [exec1, exec2]
17
18     exec1.start()
19     exec2.start()
20
21     for exec in execs:
22         exec.join()

```

```

1  khanh process: -1
2  khanh process: -2
3  khanh process: -3
4  khanh process: -4
5  ai process: -1
6  ai process: -2
7  ai process: -3
8  ai process: -4

```

Ta nhận thấy dữ liệu là không được chia sẻ giữa 2 processes vì process thứ hai đổi lại dấu của process thứ nhất sẽ khiến các phần tử của 2 processes này đảo dấu. Tiếp theo nếu chúng ta sử dụng Array trong multiprocessing thì sao ?

Top

```

1  from multiprocessing import Process, Value, Array, Lock
2  import time
3
4  def _counter_arr(arrs, process_name):
5      lock.acquire()
6      for i, el in enumerate(arrs):
7          time.sleep(0.01)
8          arrs[i] = -arrs[i]
9          print("{}: {}".format(process_name, arrs[i]))
10     lock.release()
11
12     arrs = Array('i', range(1, 5, 1))
13     lock = Lock()
14     exec1 = Process(target=_counter_arr, args=(arrs, "khanh process")) #
15     exec2 = Process(target=_counter_arr, args=(arrs, "ai process"))
16     execs = [exec1, exec2]
17
18     exec1.start()
19     exec2.start()
20
21     for exec in execs:
22         exec.join()

```

```

1  khanh process: -1
2  khanh process: -2
3  khanh process: -3
4  khanh process: -4
5  ai process: 1
6  ai process: 2
7  ai process: 3
8  ai process: 4

```

Các bạn đã thấy gì chưa ? Dữ liệu đã được chia sẻ qua lại giữa hai processes. Vậy thì chúng ta sẽ thường sử dụng shared memory khi nào ? Giả định bạn đang có một pipeline biến đổi dữ liệu gồm nhiều step khác nhau, mỗi một process sẽ phụ trách một step trong pipeline. Khi đó dữ liệu cần được shared chung giữa các process.

3.2. Shared object giữa các Process

Queue là một định dạng stack an toàn khi làm việc với multi thread và process. Chúng ta có thể tạo ra một queue và cho phép các thread, process truy cập dữ liệu mà không bị hiện tượng concurrency vì dữ liệu được truy xuất và sử dụng một lần bởi một thread hoặc process.

Bên dưới chúng ta sẽ lấy ví dụ về việc sử dụng 2 process để đọc các dữ liệu trong một queue. Hai process này tới phiên của mình sẽ lấy ra các phần tử nằm trong queue theo kiểu FIFO (First Come First Out).

Top

```

1  from multiprocessing import Process, Queue
2  import time
3
4  def _counter_queue(queue, process_name, max_count):
5      # lock.acquire()
6      while max_count:
7          time.sleep(0.01)
8          value = queue.get()
9          print("{}: {}".format(process_name, value))
10         max_count -= 1
11         # lock.release()
12
13     q = Queue()
14     for i in range(10):
15         q.put(i)
16     max_count = 5
17     # lock = Lock()
18     exec1 = Process(target=_counter_queue, args=(q, "khanh process", 5)) ;
19     exec2 = Process(target=_counter_queue, args=(q, "ai process", 5))
20     execs = [exec1, exec2]
21
22     exec1.start()
23     exec2.start()
24
25     for exec in execs:
26         exec.join()

```

```

1  khanh process: 0
2  ai process: 1
3  khanh process: 2
4  ai process: 3
5  khanh process: 4
6  ai process: 5
7  khanh process: 6
8  ai process: 7
9  khanh process: 8
10 ai process: 9

```

Như vậy không có bất kỳ một data nào được sử dụng chung giữa 2 processes nên tránh được concurrency.

3.3. Pool trong multiprocessing

Trong python chúng ta có thể sử dụng pool để tận dụng được các tính toán song song trên nhiều process một lúc. Cơ chế của pool đã loại bỏ hạn chế của GIL trong python, cho phép nhiều luồng hoạt động đồng thời và giúp đẩy nhanh quá trình tính toán.

Trong Pool chúng ta có thể khai báo nhiều workers cùng thực hiện chương trình. Các chương trình có thể thực hiện một cách bất đồng bộ thông qua hàm `apply_async()`. Tức là cho phép thực hiện song song nhiều method trên các workers. Đồng thời `apply_async()` cũng cho phép đưa vào các hàm callback để xử lý giữa liệu sau cùng.

Top

Ví dụ bên dưới chúng ta sẽ sử dụng 5 workers để tính toán bất đồng bộ bình phương của các số trong phạm vi 20. Kết quả sau khi tính sẽ được lưu vào một list.

```

1  import multiprocessing as mp
2  import time
3
4  def _square(x):
5      return x*x
6
7  def log_result(result):
8      # Hàm được gọi bất kỳ khi nào _square(i) trả ra kết quả.
9      # result_list được thực hiện trên main process, không phải pool worker
10     result_list.append(result)
11
12     def apply_async_with_callback():
13         pool = mp.Pool(processes=5)
14         for i in range(20):
15             pool.apply_async(_square, args = (i, ), callback = log_result)
16         pool.close()
17         pool.join()
18         print(result_list)
19
20     if __name__ == '__main__':
21         result_list = []
22         apply_async_with_callback()

```

```

1  [0, 1, 9, 25, 16, 36, 49, 64, 81, 100, 144, 4, 169, 225, 256, 289, 324,

```

Ta thấy thứ tự của list không theo tuần tự từ thấp tới cao do hàm được gọi bất đồng bộ.

Bên cạnh cách khởi tạo Pool cho process như trên, chúng ta còn có thể khởi tạo từ `concurrent.futures` như phần tiếp theo mình giới thiệu.

4. Process Pool vs Thread Pool

4.1. Process Pool

Trong python thì bắt đầu từ version 3.2 chúng ta có thể sử dụng module `concurrent.futures` để xử lý bất đồng bộ các tasks. Đây là một abstract layer được kế thừa trên cả hai modules là `threading` và `multiprocessing` để tạo ra một interface cho phép khởi tạo các task sử dụng pool của các processes và threads.

Để khởi tạo một Process Pool, chúng ta sử dụng `ProcessPoolExecutor` trong `concurrent.futures` module.

Top

```

1  from concurrent.futures import ProcessPoolExecutor
2  from time import sleep
3  import timeit
4
5  def _counter(counter, task_name):
6      print("Start process {}".format(task_name))
7      while (counter):
8          print("{} : {}".format(task_name, counter))
9          counter -= 1
10     print("End process {}".format(task_name))
11     return "Completed {}".format(task_name)
12
13 def _submit_process():
14     executor = ProcessPoolExecutor(max_workers=5)
15     future = executor.submit(_counter, 10, "task1")
16     print('State of future: ', future.done())
17     print('futre result: ', future.result())
18     print('State of future: ', future.done())
19
20 _submit_process()

```

```

1  Start process task1!
2  task1 : 10
3  task1 : 9
4  task1 : 8
5  task1 : 7
6  task1 : 6
7  task1 : 5
8  task1 : 4
9  task1 : 3
10 task1 : 2
11 task1 : 1
12 End process task1!
13 State of future:  False
14 futre result:  Completed task1!

```

Trong `ProcessPoolExecutor()` chúng ta cần truyền vào số lượng các worker để chạy process. Số lượng worker càng lớn thì càng nhiều threads được sinh ra để tính toán process.

Hàm `submit()`

Hàm `submit()` được sử dụng để load các task vào process pool. Tham số truyền vào là tên hàm và các đối số của hàm. Hàm `done()` để kiểm tra trạng thái của task. Lúc đầu ngay sau khi submit thì task chưa hoàn thành nên `done()` là `False`. Hàm `result()` thường được dùng để kiểm tra kết quả sau khi task cuối cùng trong process pool đã thực thi xong. Do đó trạng thái `done()` sau khi result được trả về là `True`.

Hàm `map()`

Nhắc đến hàm map trong python, nếu bạn đã có kinh nghiệm thì sẽ hiểu ngay nó sẽ map các đối số từ một list vào hàm.

Ví dụ: Để tính diện tích của các bounding box dựa trên tọa độ $(x1, x2, y1, y2)$ thì chúng ta thực hiện hàm map trong process pool như sau:

Top

```

1  from concurrent.futures import ProcessPoolExecutor
2  from concurrent.futures import as_completed
3  x1s = [5, 10, 20, 35]
4  x2s = [15, 20, 30, 55]
5  y1s = [5, 10, 10, 15]
6  y2s = [15, 20, 20, 35]
7
8  def _bbox(x1, x2, y1, y2):
9      w = x2-x1
10     h = y2-y1
11     area = w*h
12     return area
13
14  with ProcessPoolExecutor(max_workers = 5) as executor:
15      results = executor.map(_bbox, x1s, x2s, y1s, y2s)
16
17  for result in results:
18      print(result)

```



```

1      100
2      100
3      100
4      400

```

4.2. Thread Pool

Thread pool cũng tương tự như Process Pool nhưng là tập hợp của các threads thay vì processes. Các khối tạo `ThreadPoolExecutor` trên `concurrent.futures` cũng hoàn toàn tương tự như `ProcessPoolExecutor`. Ta thực hiện như sau:

```

1  from concurrent.futures import ThreadPoolExecutor
2  from time import sleep
3
4  def _counter(counter, task_name):
5      print("Start process {}".format(task_name))
6      while (counter):
7          print("{} : {}".format(task_name, counter))
8          counter -= 1
9      print("End process {}".format(task_name))
10     return "Completed {}".format(task_name)
11
12  def _submit_thread():
13     executor = ThreadPoolExecutor(max_workers=5)
14     future = executor.submit(_counter, 10, "task1")
15     print('State of future: ', future.done())
16     print('future result: ', future.result())
17     print('State of future: ', future.done())
18
19  _submit_thread()

```

Top

```

1      Start process task1!
2      task1 : 10
3      task1 : 9
4      task1 : 8
5      task1 : 7
6      task1 : 6
7      task1 : 5
8      task1 : 4
9      task1 : 3
10     task1 : 2
11     task1 : 1
12     End process task1!
13     State of future: True
14     futre result: Completed task1!

```

So sánh thời gian xử lý của process và thread. Để đo lường, các bạn cần cài package cProfile .

```

1      import cProfile
2      cProfile.run('_submit_process()')
3      cProfile.run('_submit_thread()')

1      State of future: True
2          1268 function calls in 0.072 seconds
3
4      ...
5
6      State of future: True
7          249 function calls in 0.001 seconds
8      ...

```

Với cùng các tác vụ như nhau thì ta thấy thời gian thực thi của Thread Pool chỉ là 0.001 seconds, nhanh thời gian thực thi của Process Pool là 0.072 seconds. Lý do là vì thread là một phiên bản light weight hơn process rất nhiều. Bạn có thể nhận thấy điều này một cách trực quan thông qua số hàm được gọi ở cả hai phương pháp.

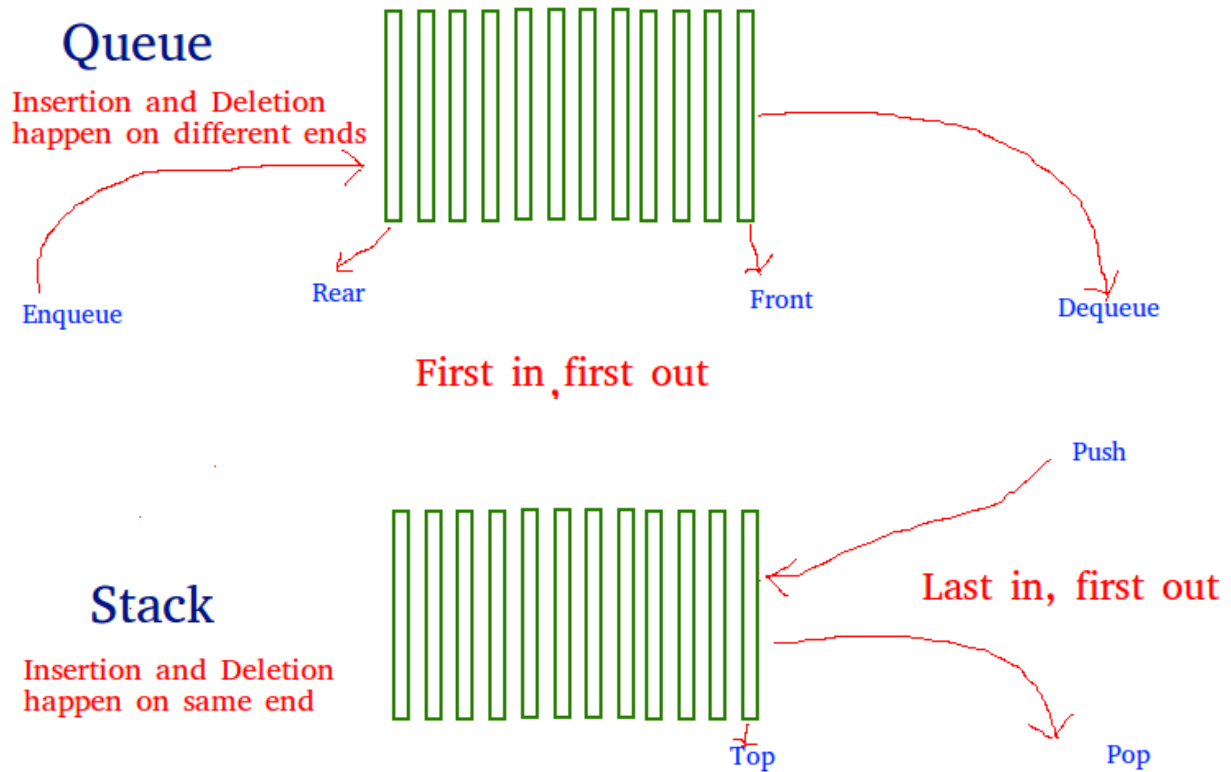
Vậy lựa chọn thế nào giữa process pool và thread pool?

Chúng ta đã biết rằng khi sử dụng threads thì sẽ có lợi về I/O vì các threads có thể chia sẻ data qua lại lẫn nhau. Còn giữa các processes thì data được sử dụng hoàn toàn độc lập nên không có lợi về I/O. Tuy nhiên khi sử dụng process thì chúng ta sẽ được allocate về CPU, Memomory,... nên lời khuyên là nếu task của bạn gặp phải giới hạn về I/O bound thì nên sử dụng thread pool và giới hạn về CPUs bound thì nên sử dụng process pool.

5. Queue

Khi chạy một process trên multiple-threads, queue thường được sử dụng để làm nơi lưu trữ dữ liệu chung giữa các threads với nhau.

Top



Hình 3: Cấu trúc FIFO và LIFO của một queue.

Queue là một cấu trúc dữ liệu tuyến tính (linear data structure). Nó có tính chất tương tự như list. Cho phép chúng ta thêm, sửa, xóa, truy xuất các phần tử bên trong. Trong python, Queue có ưu điểm lớn hơn list đó là tốc độ truy xuất nhanh hơn. Độ phức tạp thời gian (time complexity) của queue là $O(1)$ trong khi của list là $O(n)$. Queue là một lựa chọn thay thế tốt hơn cho list trong trường hợp dữ liệu của bạn có số lượng phần tử lớn.

Khi làm việc với Queue bạn có thể truy xuất các phần tử bên trong khối theo kiểu FIFO (first in first out) hoặc LIFO (last in first out) thông qua hàm `pop()` .

Queue thường được sử dụng trong các tác vụ liên quan tới threads synchronous. Các thread sẽ sử dụng chung một dữ liệu và thay đổi các phần tử bên trong nó một cách tuần tự.

```

1  from concurrent.futures import ThreadPoolExecutor
2  import queue
3
4
5  def _sum_queue(name, work_queue):
6      sum = 0
7      while not work_queue.empty():
8          print(f"Task {name} running")
9          count = work_queue.get()
10         sum += count
11         print(f"Task {name} total: {sum}")
12         return sum
13
14  def task(name, work_queue):
15      if work_queue.empty():
16          print(f"Task {name} nothing to do")
17      else:
18          print("Start ThreadPoolExecutor!")
19          with ThreadPoolExecutor(max_workers = 5) as executor:
20              print("Submit task!")
21              future = executor.submit(_sum_queue, name, work_queue)
22              sum = future.result()
23              return sum
24
25  # Create the queue of work
26  work_queue = queue.Queue()
27
28  # Put some work in the queue
29  for work in [15, 10, 5, 2]:
30      work_queue.put(work)
31
32  # Create some synchronous tasks
33  tasks = [("one", work_queue), ("two", work_queue)]
34
35  # Run the tasks
36  for n, q in tasks:
37      print(task(n, q))

```

```

1  Start ThreadPoolExecutor!
2  Submit task!
3  Task one running
4  Task one running
5  Task one running
6  Task one running
7  Task one total: 32
8  32
9  Task two nothing to do
10 None

```

Trong ví dụ trên giải sử chúng ta có hai threads hoạt động một cách synchronous là one và two . Hai threads này sử dụng chung một nguồn dữ liệu là work_queue . Khi thread one chạy xong thì toàn bộ các phần tử của queue đã được trích xuất xong nên ở thread two chúng ta không có gì để chạy tiếp.

Top

6. Kết luận

Như vậy qua bài viết này mình đã giới thiệu với các bạn những khái niệm cơ bản nhất trong xử lý song song như thread, process, khác biệt giữa chúng và các khái niệm về đồng bộ, bất đồng bộ, concurrency. Thông qua ví dụ minh họa chi tiết và cách khởi tạo và áp dụng thread và process trong một chương trình bạn đọc sẽ có thể áp dụng vào các project của mình.

Đồng thời để loại bỏ hạn chế của GIL trong python, tận dụng được tính toán song song trên nhiều CPUs thì bạn đọc có thể áp dụng Process Pool, Thread Pool như một giải pháp hữu hiệu.

Để xây dựng bài viết này mình đã khảo cứu và tổng hợp từ rất nhiều nguồn dữ liệu bên dưới.

7. Tài liệu

<https://docs.python.org/3.1/library/multiprocessing.html>

<https://docs.python.org/3/library/queue.html#queue.Queue>

<https://www.geeksforgeeks.org/stack-in-python/?ref=lbp>

<https://realpython.com/python-async-features/>

https://www.tutorialspoint.com/concurrency_in_python/concurrency_in_python_pool_of_processes.htm

<https://medium.com/@bfortuner/python-multithreading-vs-multiprocessing-73072ce5600b>

<https://docs.python.org/3/library/asyncio-subprocess.html>

<https://realpython.com/python-async-features/>

Top