

Turbo Python Performance To Achieve 100x Faster

June 2023

Challenges and Perceptions:

- Python GIL design is slow**
- Python is a slow at runtime**
- Python performance is slower comparing to c++ and java**

Can it be improved?

- Let's find it out**

Techniques used in this performance testing

- Step up and use built-in functions**
- Use vectorization**
- Use math functions**
- Use multi-processing, concurrency**
- Memoization and caching**
- Use different language at server backend**
- Engineering thoughts**

Note:

**Often many things can impact python runtime performance
From hardware, cpu, memory, latency in addition to code
Ideas presented here focus on python source code only**

```
In [1]: 1 # Use built-in functions and libraries, they are tested and optimized
        2 import string
        3 def upper_basic(n):
        4     newList = []
        5     for w in string.ascii_lowercase*n:
        6         newList.append(w.upper())
```

```
In [2]: 1 %timeit upper_basic(1000)
```

4.12 ms \pm 23.1 μ s per loop (mean \pm std. dev. of 7 runs, 100 loops each)

```
In [3]: 1 def upper_o2(n):
        2     newList = map(str.upper, string.ascii_lowercase*n)
```

```
In [72]: 1 %timeit upper_o2(1000)
```

1.22 μ s \pm 18.6 ns per loop (mean \pm std. dev. of 7 runs, 1,000,000 loops each)

60% better when using map() which does elementwise operation

Use vectorization

apply operations to all elements of an array in one go

"for" loop manipulates one row at a time

```
In [20]: 1 def find_sum(n):  
2         total = 0  
3         for i in range(n):  
4             total += i
```

```
In [21]: 1 %timeit find_sum(1_000_000)
```

74.9 ms \pm 1.22 ms per loop (mean \pm std. dev. of 7 runs, 10 loops each)

```
In [22]: 1 import numpy as np  
2         def find_sum_vector(n):  
3             total = 0  
4             total = np.sum(np.arange(n))
```

```
In [23]: 1 %timeit find_sum_vector(1_000_000)
```

2.08 ms \pm 78.9 μ s per loop (mean \pm std. dev. of 7 runs, 100 loops each)

```
In [113]: 1 (74.9-2.08)/74.9
```

```
Out[113]: 0.9722296395193591
```

97% better when using vectorization in numpy

Deep learning multi-linear regression calculations

$$y = m_1x_1 + m_2x_2 + m_3x_3 + m_4x_4 + m_5x_1 + c$$

Use loop for million of rows of calculations is slow

Vectorization is the optimal solution

```
In [38]: 1 # create random data
2 import numpy as np
3 m = np.random.rand(1,5)
4 n = np.random.rand(100000,5)
5 m.shape, n.shape
6
```

```
Out[38]: ((1, 5), (100000, 5))
```

```
In [39]: 1 # use loop for calculations
2 import numpy as np
3
4 def loop_reg_sum(col, row):
5     m = np.random.rand(1,col)
6     n = np.random.rand(row,col)
7     result = []
8     for i in range(row):
9         total = 0
10        for j in range(col):
11            total += n[i][j]*m[0][j]
12        # print(i, total)
13        result.append(total)
14
```

```
In [40]: 1 %timeit loop_reg_sum(5, 100_000)
```

407 ms ± 7.12 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```
In [41]: 1 # use vectorization
2 def vec_reg_sum(col, row):
3     m = np.random.rand(1,col)
4     n = np.random.rand(row,col)
5     result = np.dot(n, m.T)
6
```

```
In [42]: 1 %timeit vec_reg_sum(5, 100_000)
```

6.45 ms ± 270 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)

407 vs 6.45 Regression using python loop and numpy vectorization

```

In [49]: 1 # use built-in sum
          2 def sum_range(n=1_000_000):
          3     return sum(range(n))

In [50]: 1 # use numpy (implemented in c, faster)
          2 import numpy
          3 def sum_numpy(n=1_000_000):
          4     return numpy.sum(numpy.arange(n))    # this is a one c call, but a whole array is created in memory
          5

In [51]: 1 # use math knowledge
          2 def sum_math(n=1_000_000):
          3     return (n * (n-1)) // 2

In [52]: 1 import timeit
          2
          3 print('while loop\t\t', timeit.timeit(while_loop, number = 1))
          4 print('for loop\t\t', timeit.timeit(for_loop, number = 1))
          5 print('for_loop_with_increment\t\t', timeit.timeit(for_loop_with_increment, number = 1))
          6 print('for_loop_with_test\t\t', timeit.timeit(for_loop_with_test, number = 1))
          7 print('for_loop_with_increment_and_test\t\t', timeit.timeit(for_loop_with_increment_and_test, number = 1))
          8 print('sum_range\t\t', timeit.timeit(sum_range, number = 1))
          9 print('sum_numpy\t\t', timeit.timeit(sum_numpy, number = 1))
         10 print('sum_math\t\t', timeit.timeit(sum_math, number = 1))
         11
         12 # python programming consideration
         13 # use math formula
         14 # use c implementation
         15 # use built-in function, sum, map ... which loops for you
         16 # for or while loop
         17
         18
while loop          0.11545043099999993
for loop            3.713000069183181e-06
for_loop_with_increment    0.11514549399998941
for_loop_with_test        0.0975002409999206
for_loop_with_increment_and_test    0.1432448260000001
sum_range           0.022303372999999738
sum_numpy           0.005475752999927863
sum_math            2.3169999394667684e-06

```

Different loops and their performances

```
In [53]: 1 ## memoization or cache to optimize
2 # useful for recursive functions, or operations used over and over again
3 # you don't want to repeat to calculate values again
```

```
In [115]: 1 # use cache dict
2 from time import perf_counter
3 from functools import wraps
4
5 def memoize(func):
6     cache = {}
7
8     @wraps(func)
9     def wrapper(*args, **kwargs):
10         key = str(args) + str(kwargs)
11         if key not in cache:
12             cache[key] = func(*args, **kwargs)
13         return cache[key]
14     return wrapper
15
```

```
In [116]: 1 # fibonacci using memoize
2 def fibonacci_plain(n=100) -> int:
3     if n < 2:
4         return n
5     return fibonacci_plain(n-1) + fibonacci_plain(n-2)
```

```
In [*]: 1 # no memoization call, very slow, cpu humming, 20 mins still running, killed this cell
2 start = perf_counter()
3 fibonacci_plain()
4 end = perf_counter()
5 print(end-start)
```

```
In [57]: 1 # fibonacci using memoize
2 @memoize
3 def fibonacci(n=1000) -> int:
4     if n < 2:
5         return n
6     return fibonacci(n-1) + fibonacci(n-2)
```

```
In [58]: 1 # get result instantly
2 print('fibonacci with memorize\t\t', timeit.timeit(fibonacci, number = 1))
```

fibonacci with memorize 0.0036856550000265997

Memoization and caching reducing intermediate operations

```

27 outputs = []
28 for url in urls.values():
29     print(url)
30     outputs = outputs + [requests.get(url).text]
31     #print(outputs)
32
33 count_https = []
34 count_http = []
35 for output in outputs:
36     count_https += re.findall("https://", output)
37     count_http += re.findall("http://", output)
38
39 print(len(count_https), len(count_http))
40
41 # index = 0
42 # while count_https[index]:
43 #     if index >= len(count_https):
44 #         break
45 #     index += 1
46
47 start = time.perf_counter()
48 count_words_in_web_page()
49 elapsed = time.perf_counter() - start
50 print(f'{elapsed:.2f} seconds')
51

```

```

https://google.com
https://yahoo.com
https://microsoft.com
https://google.com
https://apple.com
https://ibm.com
https://amazon.com
https://twitter.com
https://tiktok.com
https://oracle.com
https://intel.com
https://tesla.com
https://nasa.com
https://ebay.com
https://wikipedia.com
3071 732
10.50 seconds

```

```

14 "2": "https://yahoo.com",
15 "3": "https://microsoft.com",
16 "4": "https://google.com",
17 "5": "https://apple.com",
18 "6": "https://ibm.com",
19 "7": "https://amazon.com",
20 "8": "https://twitter.com",
21 "9": "https://tiktok.com",
22 "10": "https://oracle.com",
23 "11": "https://intel.com",
24 "12": "https://tesla.com",
25 "13": "https://nasa.com",
26 "14": "https://ebay.com",
27 "15": "https://wikipedia.com"
28 }
29
30 # mark as async
31 async def count_words_in_web_page_async():
32     outputs = []
33
34     async with httpx.AsyncClient() as client:
35         tasks = (client.get(url) for url in urls.values())
36         reqs = await asyncio.gather(*tasks) # waits for task, but await till all donee
37
38         outputs = [req.text for req in reqs]
39         #print(outputs)
40
41     count_https, count_http = [], []
42     for output in outputs:
43         count_https += re.findall("https://", output) # text processing, not use pre-compiled re
44         count_http += re.findall("http://", output)
45     # print(count_https)
46     # print(count_http)
47
48 start = time.perf_counter()
49 await count_words_in_web_page_async() # schedule func to run
50 # asyncio.run(count_words_in_web_page_async()) # for python>3.7 and ipython < 7.0
51 elapsed = time.perf_counter() - start
52 print(f'{elapsed:.2f} seconds')
53
54
55

```

1.03 seconds

10.50 vs 1.03 - Use async for web text scraping


```

(base) user-2:bin user$ cat main.rs
extern crate webserver;

use webserver::ThreadPool;
use std::net::TcpListener;
use std::io::prelude::*;
use std::net::TcpStream;
use std::fs::File;
use std::thread;
use std::time::Duration;

fn main() {

    let listener = TcpListener::bind("127.0.0.1:7878").unwrap();
    let pool = ThreadPool::new(8);

    for stream in listener.incoming() {

        let stream = stream.unwrap();

        pool.execute(|| {
            handle_connection(stream);
        });

    }

}

fn handle_connection(mut stream: TcpStream) {

    let mut buffer = [0; 512];

    stream.read(&mut buffer).unwrap();

    let get = b"GET / HTTP/1.1\r\n";
    let sleep = b"GET /sleep HTTP/1.1\r\n";

    let (status_line, filename) = if buffer.starts_with(get) {
        ("HTTP/1.1 200 OK\r\n\r\n", "hello.html")
    } else if buffer.starts_with(sleep) {
        thread::sleep(Duration::from_secs(5));
        ("HTTP/1.1 200 OK\r\n\r\n", "hello.html")
    } else {
        ("HTTP/1.1 404 NOT FOUND\r\n\r\n", "404.html")
    };

    let mut file = File::open(filename).unwrap();
    let mut contents = String::new();

    file.read_to_string(&mut contents).unwrap();

    let response = format!("{}", status_line, contents);

    stream.write(response.as_bytes()).unwrap();
    stream.flush().unwrap();

}

```

Rust multi thread server

Python web api call

```
url = 'http://localhost:7878/'

def fetch(session, url):
    with session.get(url) as response:
        #print(response)
        pass

@timer(1, 1)
def main():
    with requests.Session() as session:
        for _ in range(5000):
            fetch(session, url)

import requests
import timeit
from multiprocessing.pool import Pool

url = 'http://localhost:7878/'

def fetch(session, url):
    with session.get(url) as response:
        #print(response)
        pass

def timer(number, repeat):
    def wrapper(func):
        runs = timeit.repeat(func, number=number, repeat=repeat)
        print(sum(runs) / len(runs))
    return wrapper

if __name__ == "__main__":

    @timer(1, 1)
    def task():
        with Pool() as pool:
            with requests.Session() as session:
                pool.starmap(fetch, [(session, url) for _ in range(5000)])

(base) user-2:client user$
(base) user-2:client user$ cat ../readme.txt
# start server
(base) user-2:rust_web_server_concurrent user$ cargo run
    Finished dev [unoptimized + debuginfo] target(s) in 0.00s
    Running `target/debug/main`

# start 01_simple-http.py
# synchronous calls
(base) user-2:client user$ python 01_simple_http_sync.py
11.093317224

# multiprocessing using multi-cores to run multi processes
(base) user-2:client user$ python 02_multi_processing_http.py
2.7904895300000003
```

Python web api call

```
url = 'http://localhost:7878/'

def fetch(session, url):
    with session.get(url) as response:
        #print(response)
        pass

@timer(1, 1)
def main():
    with requests.Session() as session:
        for _ in range(5000):
            fetch(session, url)

import requests
import timeit
from multiprocessing.pool import Pool

url = 'http://localhost:7878/'

def fetch(session, url):
    with session.get(url) as response:
        #print(response)
        pass

def timer(number, repeat):
    def wrapper(func):
        runs = timeit.repeat(func, number=number, repeat=repeat)
        print(sum(runs) / len(runs))
    return wrapper

if __name__ == "__main__":

    @timer(1, 1)
    def task():
        with Pool() as pool:
            with requests.Session() as session:
                pool.starmap(fetch, [(session, url) for _ in range(5000)])

(base) user-2:client user$
(base) user-2:client user$ cat ../readme.txt
# start server
(base) user-2:rust_web_server_concurrent user$ cargo run
    Finished dev [unoptimized + debuginfo] target(s) in 0.00s
    Running `target/debug/main`

# start 01_simple-http.py
# synchronous calls
(base) user-2:client user$ python 01_simple_http_sync.py
11.093317224

# multiprocessing using multi-cores to run multi processes
(base) user-2:client user$ python 02_multi_processing_http.py
2.7904895300000003
```

```
In [2]: 1 !python -V
```

Python 3.10.11

```
In [1]: 1 import math
        2 import numpy as np
```

Factorial goes faster

```
In [39]: 1 # basic factorial
        2 %timeit math.prod(range(1, 150))
```

10.6 μ s \pm 32.3 ns per loop (mean \pm std. dev. of 7 runs, 100,000 loops each)

```
In [33]: 1 math.prod(range(1,6))
```

Out[33]: 120

```
In [34]: 1 def f(x):
        2     return x * f(x-1) if x > 1 else 1
```

```
In [28]: 1 f(5)
```

Out[28]: 120

```
In [40]: 1 %timeit math.factorial(150)
```

1.8 μ s \pm 3.93 ns per loop (mean \pm std. dev. of 7 runs, 1,000,000 loops each)

```
In [44]: 1 (1.8/10.6)
```

Out[44]: 0.169811320754717

Use new versions

compute factorials like a boss

Winning ideas:

Left shift is cheaper than multiplying by two

Pulling out events leaves recurring odd factories

Dynamic programming reuses previously computed odd factorials

```
In [75]: 1 s = ""
2 so = ""
3 se = ""
4 for i in range(1,20):
5     s += str(i) + '*'
6     if i%2 == 0:
7         se += str(i) + '*'
8     else:
9         so += str(i) + '*'
10    print(s, so, se)

1*2*3*4*5*6*7*8*9*10*11*12*13*14*15*16*17*18*19* 1*3*5*7*9*11*13*15*17*19* 2*4*6*8*10*12*14*16*18*
```

```
In [71]: 1 1*2*3*4*5*6*7*8*9*10*11*12*13*14*15*16*17*18*19
```

Out[71]: 121645100408832000

```
In [79]: 1 # collect even terms, next divide each even by 2
2 1*3*5*7*9*11*13*15*17*19 *2*4*6*8*10*12*14*16*18
```

Out[79]: 121645100408832000

```
In [81]: 1 # divide terms by two and replace with left shift
2 1*3*5*7*9*11*13*15*17*19 *1*2*3*4*5*6*7*8*9 << 9
```

Out[81]: 121645100408832000

```
In [84]: 1 1*3*5*7*9*11*13*15*17*19 *1*3*5*7*9 * 1*2*3*4 <<13
```

Out[84]: 121645100408832000

```
In [86]: 1 1*3*5*7*9*11*13*15*17*19 *1*3*5*7*9 * 1*3 *1*2 <<15
```

Out[86]: 121645100408832000

```
In [88]: 1 # replace even term with left shift
2 1*3*5*7*9*11*13*15*17*19 *1*3*5*7*9 * 1*3 <<16
```

Out[88]: 121645100408832000

```
In [93]: 1 # factor-out common subsequences and replace with powers
2 (1*3)**3 *(5*7*9)**2 *(11*13*15*17*19)**1 <<16
```

Out[93]: 121645100408832000

How python 3 implement factorial() performance

Faster CPython

performance

1/8

CPython 3.11 is an average of **25% faster** than CPython 3.10 as measured with the **pyperformance** benchmark suite, when compiled with GCC on Ubuntu Linux. Depending on your workload, the overall speedup could be 10–60%.

This project focuses on two major areas in Python: **Faster Startup** and **Faster Runtime**. Optimizations not covered by this project are listed separately under **Optimizations**.

Faster Startup

Frozen imports / Static code objects

Python caches **bytecode** in the `__pycache__` directory to speed up module loading.

Previously in 3.10, Python module execution looked like this:

```
Read __pycache__ -> Unmarshal -> Heap allocated code object -> Evaluate
```

In Python 3.11, the core modules essential for Python startup are “frozen”. This means that their **Code Objects** (and bytecode) are statically allocated by the interpreter. This reduces the steps in module execution process to:

```
Statically allocated code object -> Evaluate
```

Interpreter startup is now 10–15% faster in Python 3.11. This has a big impact for short-running programs using Python.

(Contributed by Eric Snow, Guido van Rossum and Kumar Aditya in many issues.)

Faster Runtime

Cheaper, lazy Python frames

Python frames, holding execution information, are created whenever Python calls a Python function. The following are new frame optimizations:

- Streamlined the frame creation process.
- Avoided memory allocation by generously re-using frame space on the C stack.
- Streamlined the internal frame struct to contain only essential information. Frames previously held extra debugging and memory management information.

Faster python 3.11

Conclusions

- Python is NOT the programming language to solve all problems**
- BUT Python is a popular language for many use cases**
- Performance can be improved**
- Python is improving, by itself and community**

- In the era of ML/AI, no doubt python will become popular**
- To learn, experiment, build quickly MVP in less time**

Q & A