# 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 optimzied
           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 ± 23.1 µs per loop (mean ± 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 µs ± 18.6 ns per loop (mean ± 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 ± 1.22 ms per loop (mean ± 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 ± 78.9 µs per loop (mean ± 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**

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + m_4 x_4 + m_5 x_1 + c$$

Use loop for million of rows of calculations is slow
Vectorization is the optimal solution

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

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

In [39]:
```python
# use loop for calculations
import numpy as np

def loop_reg_sum(col, row):
    m = np.random.rand(1,col)
    n = np.random.rand(row,col)
    result = []
    for i in range(row):
        total = 0
        for j in range(col):
            total += n[j][j]*m[0][j]
#           print(i, total)
        result.append(total)
```

In [40]:
```python
%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]:
```python
# use vectorization
def vec_reg_sum(col, row):
    m = np.random.rand(1,col)
    n = np.random.rand(row,col)
    result = np.dot(n, m.T)
```

In [42]:
```python
%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               0.1154504309999993
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 memoize\t\t', timeit.timeit(fibonacci, number = 1))

            fibonacci with memoize          0.0036856550000265997
```

# Memoization and caching reducing intermediate operations

```python
27    outputs = []
28    for url in urls.values():
29        print(url)
30        outputs = outputs + [requests.get(url).tex
31        #print(outputs)
32
33    count_https = []
34    count_http = []
35    for output in outputs:
36        count_https += re.findall("https://", outpu
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
```

```python
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
49 start = time.perf_counter()
50 await (count_words_in_web_page_async())    # schedule func to run
51 # asyncio.run(count_words_in_web_page_async())   # for python>3.7 and ipython < 7.0
52 elapsed = time.perf_counter() - start
53 print(f'{elapsed:.2f} seconds')
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

```
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 processeso
(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 processeso
(base) user-2:client user$ python 02_multi_processing_http.py
2.7904895300000003
```

Python web api call

```
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 µs ± 32.3 ns per loop (mean ± 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 µs ± 3.93 ns per loop (mean ± 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]:
```python
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]:
```python
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]:
```python
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]:
```python
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]:
```python
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]:
```python
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]:
```python
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]:
```python
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

## Conclusions

-Python is the programming language to solve all problems
-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**