

Ravella, Kiran

09/17/2019

Harvard University Extension - Principles of Big Data Processing e88

Homework 2: Vertical and Horizontal Scaling , Shared State Management

This document is a template for your solutions submission. You are free to add additional information in this submission if you would like. Extra screenshots and extra documentation are appreciated. Screenshots must always be viewable. If a screenshot is too blurry or chopped off in a key area you will not receive full credit for it.

Make sure to also submit all your source code (.java files , .py files or whatever language you are using) - in a separate archive, named <LastName>_<FirstName>_HW2.zip

Please identify which problems were completed. If any were incomplete, please identify where you encountered problems.

Problem 1: 100% complete
Problem 2: 100% complete
Problem 3: 100% complete
Problem 4: 100% complete
Problem 5 Bonus: 0% Complete

Problem 1: CPU Analysis [points: 25]

Paste your source code into the following area [10 points]

```
# This program creates a CPU-Intensive workload.
import argparse
import os
import time
import multiprocessing

#Using argparse library to parse the input arguments.
prog = "hw2"
desc = "Starts a specified number of threads that do CPU-Intensive work"
parser = argparse.ArgumentParser(prog=prog, description=desc)
parser.add_argument('--numThreads', '-n', default=4, type=int)

parsed_args = parser.parse_args()
numThreads = parsed_args.numThreads

print("Number of Threads: ", numThreads)

#Function that calculates the Fibonacci series
def Fibonacci(n):
    if n==1:
        return 0
    elif n==2:
        return 1
    else:
        return Fibonacci(n-1)+Fibonacci(n-2)

#Function that loops infinitely and puts load on the CPU.
```

```
def create_cpu_load(thread_num):
    while(True):
        time.sleep(1)
        print("In Thread :", thread_num, ", Process ID: ", os.getpid())
        Fibonacci(100)
        print("In Thread:", thread_num," ,Completed Fibonacci")

#The following code creates the number of threads specified in the input argument
threads= []
for thread_count_id in range(numThreads):
    t = multiprocessing.Process(target = create_cpu_load, args=(thread_count_id,))
    threads.append(t)
    t.start()

for current_thread in threads:
    current_thread.join()
```

Provide your table or graphs demonstrating the results of running this code with 2, 4, and 16 threads on a 4 CPU machine: [5 points]

Running Code with 2 Threads on a c4.xlarge (4 CPU) instance:

```
[[centos@ip-172-31-13-96 ~]$ python cpu_intensive_multi_processing.py --numThreads 2
('Number of Threads: ', 2)
('In Thread :', 0, ', Process ID: ', 13118)
('In Thread :', 1, ', Process ID: ', 13119)
```

hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 156x29

1	[]	0.7%	Tasks: 36, 57 thr, 79 kthr; 3 running
2	[]	0.0%	Load average: 2.00 1.51 0.73
3	[]	100.0%	Uptime: 00:15:48
4	[]	100.0%	
Mem	[]	289M/7.15G	
Swp	[]	OK/OK	

PPID	CPU	TGID	PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
12559	1	13117	13117	centos	20	0	152M	7040	2268	S	0.0	0.1	0:00.01	python cpu_intensive_multi_processing.py --numThreads 2
13117	4	13119	13119	centos	20	0	152M	5368	596	R	99.7	0.1	7:01.72	python cpu_intensive_multi_processing.py --numThreads 2
13117	3	13118	13118	centos	20	0	152M	5552	740	R	99.7	0.1	7:01.72	python cpu_intensive_multi_processing.py --numThreads 2
1	4	858	858	root	20	0	549M	16832	6024	S	0.0	0.2	0:00.23	/usr/bin/python -Es /usr/sbin/tuned -l -P

Running Code with 4 Threads on a c4.xlarge (4 CPU) instance:

```
[[centos@ip-172-31-13-96 ~]$ python cpu_intensive_multi_processing.py --numThreads 4
('Number of Threads: ', 4)
('In Thread :', 0, ', Process ID: ', 13583)
('In Thread :', 1, ', Process ID: ', 13584)
('In Thread :', 2, ', Process ID: ', 13585)
('In Thread :', 3, ', Process ID: ', 13586)
```

hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 156x29

1	[]	100.0%	Tasks: 39, 57 thr, 79 kthr; 5 running
2	[]	100.0%	Load average: 4.04 3.79 2.43
3	[]	100.0%	Uptime: 00:28:01
4	[]	100.0%	
Mem	[]	295M/7.15G	
Swp	[]	OK/OK	

PPID	CPU	TGID	PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
12559	3	13582	13582	centos	20	0	152M	7044	2268	S	0.0	0.1	0:00.01	python cpu_intensive_multi_processing.py --numThreads 4
13582	3	13586	13586	centos	20	0	152M	5372	596	R	100.	0.1	9:50.79	python cpu_intensive_multi_processing.py --numThreads 4
13582	1	13585	13585	centos	20	0	152M	5372	596	R	100.	0.1	9:50.60	python cpu_intensive_multi_processing.py --numThreads 4
13582	4	13584	13584	centos	20	0	152M	5372	596	R	100.	0.1	9:50.10	python cpu_intensive_multi_processing.py --numThreads 4
13582	2	13583	13583	centos	20	0	152M	5556	740	R	99.5	0.1	9:49.02	python cpu_intensive_multi_processing.py --numThreads 4
1	4	858	858	root	20	0	549M	16832	6024	S	0.0	0.2	0:00.28	/usr/bin/python -Es /usr/sbin/tuned -l -P

Running Code with 16 Threads on a c4.xlarge (4 CPU) instance:

```
[centos@ip-172-31-13-96 ~]$ python cpu_intensive_multi_processing.py --numThreads 16
('Number of Threads: ', 16)
('In Thread :', 2, ' ', Process ID: ', 14205)
('In Thread :', 0, ' ', Process ID: ', 14203)
('In Thread :', 1, ' ', Process ID: ', 14204)
('In Thread :', 4, ' ', Process ID: ', 14207)
('In Thread :', 6, ' ', Process ID: ', 14209)
('In Thread :', 3, ' ', Process ID: ', 14206)
('In Thread :', 5, ' ', Process ID: ', 14208)
('In Thread :', 8, ' ', Process ID: ', 14211)
('In Thread :', 9, ' ', Process ID: ', 14212)
('In Thread :', 10, ' ', Process ID: ', 14213)
('In Thread :', 11, ' ', Process ID: ', 14214)
('In Thread :', 12, ' ', Process ID: ', 14215)
('In Thread :', 13, ' ', Process ID: ', 14216)
('In Thread :', 7, ' ', Process ID: ', 14210)
('In Thread :', 14, ' ', Process ID: ', 14217)
('In Thread :', 15, ' ', Process ID: ', 14218)
```

```
hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 156x32

1 [|||||] 100.0% Tasks: 50, 57 thr, 80 kthr; 17 running
2 [|||||] 100.0% Load average: 16.03 11.69 6.40
3 [|||||] 100.0% Uptime: 00:35:11
4 [|||||] 100.0%

Mem[|||||] 302M/7.15G
Swp[|||||] OK/OK

PPID CPU Tgid PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
12559 1 14202 14202 centos 20 0 152M 7044 2268 S 0.0 0.1 0:00.02 python cpu_intensive_multi_processing.py --numThreads 16
14202 4 14218 14218 centos 20 0 152M 5372 596 R 25.0 0.1 1:16.40 python cpu_intensive_multi_processing.py --numThreads 16
14202 4 14217 14217 centos 20 0 152M 5372 596 R 25.0 0.1 1:16.72 python cpu_intensive_multi_processing.py --numThreads 16
14202 4 14216 14216 centos 20 0 152M 5372 596 R 24.4 0.1 1:16.90 python cpu_intensive_multi_processing.py --numThreads 16
14202 2 14215 14215 centos 20 0 152M 5372 596 R 25.0 0.1 1:16.44 python cpu_intensive_multi_processing.py --numThreads 16
14202 3 14214 14214 centos 20 0 152M 5372 596 R 26.4 0.1 1:16.68 python cpu_intensive_multi_processing.py --numThreads 16
14202 2 14213 14213 centos 20 0 152M 5372 596 R 24.4 0.1 1:16.94 python cpu_intensive_multi_processing.py --numThreads 16
14202 2 14212 14212 centos 20 0 152M 5372 596 R 23.7 0.1 1:17.02 python cpu_intensive_multi_processing.py --numThreads 16
14202 1 14211 14211 centos 20 0 152M 5372 596 R 27.0 0.1 1:16.39 python cpu_intensive_multi_processing.py --numThreads 16
14202 3 14210 14210 centos 20 0 152M 5372 596 R 25.7 0.1 1:16.90 python cpu_intensive_multi_processing.py --numThreads 16
14202 1 14209 14209 centos 20 0 152M 5372 596 R 25.7 0.1 1:16.93 python cpu_intensive_multi_processing.py --numThreads 16
14202 3 14208 14208 centos 20 0 152M 5368 596 R 23.7 0.1 1:16.44 python cpu_intensive_multi_processing.py --numThreads 16
14202 4 14207 14207 centos 20 0 152M 5368 596 R 24.4 0.1 1:16.63 python cpu_intensive_multi_processing.py --numThreads 16
14202 1 14206 14206 centos 20 0 152M 5368 596 R 25.7 0.1 1:16.98 python cpu_intensive_multi_processing.py --numThreads 16
14202 2 14205 14205 centos 20 0 152M 5368 596 R 24.4 0.1 1:16.90 python cpu_intensive_multi_processing.py --numThreads 16
14202 1 14204 14204 centos 20 0 152M 5368 596 R 22.4 0.1 1:16.41 python cpu_intensive_multi_processing.py --numThreads 16
14202 3 14203 14203 centos 20 0 152M 5536 740 R 25.7 0.1 1:16.79 python cpu_intensive_multi_processing.py --numThreads 16
1 4 858 858 root 20 0 549M 16832 6024 S 0.0 0.2 0:00.30 /usr/bin/python -Es /usr/sbin/tuned -l -P
```

Number of Threads	Load Average over last minute
2	2
4	4.04
16	16.03

What can you summarize about the results? [3 points]

When the number of threads are less than the number of CPUs, then the program is executed in the same number of CPUs as the number of threads and those CPUs are fully utilized. The

remaining CPUs are not utilized. In our case, we ran the program with 2 threads and 2 CPUs were 100% Utilized whereas the rest were close to 0% utilization.

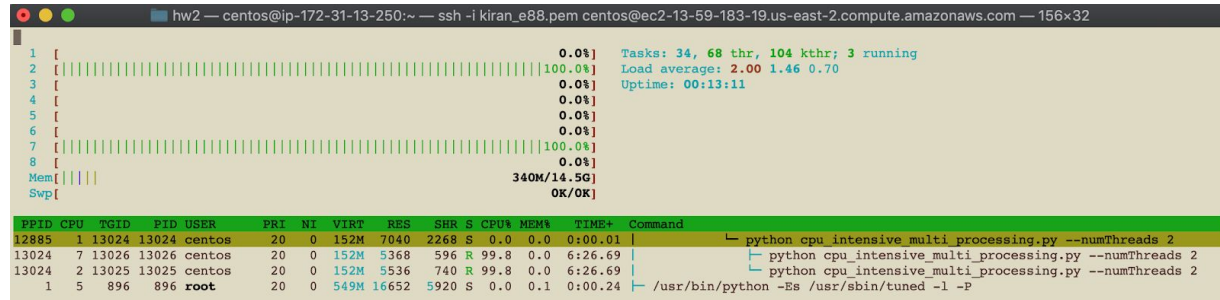
When the number of threads is equal to the number of CPUs, all the CPUs are fully utilized. In our case, we ran the program with 4 threads and all 4 CPUs were 100% utilized.

When the number of threads is more than the number of CPUs, then all the CPUs are fully utilized and the load average indicates that the CPUs are overloaded and that there are processes waiting for CPU time. In our case, we ran with 16 threads, so all 4 CPUs were 100% utilized and there were a lot of processes waiting for CPU time, which is more than 1200%.

Provide your table or graphs demonstrating the results of running this code with 2, 4, and 16 threads on an 8 CPU machine: [5 points]

Running Code with 2 Threads on a c4.2xlarge (8 CPU) instance:

```
[[centos@ip-172-31-13-250 ~]$ python cpu_intensive_multi_processing.py --numThreads 2
('Number of Threads: ', 2)
('In Thread :', 0, ', Process ID: ', 13025)
('In Thread :', 1, ', Process ID: ', 13026)
```



Running Code with 4 Threads on a c4.2xlarge (8 CPU) instance:

```
[[centos@ip-172-31-13-250 ~]$ python cpu_intensive_multi_processing.py --numThreads 4
('Number of Threads: ', 4)
('In Thread :', 0, ', Process ID: ', 13374)
('In Thread :', 2, ', Process ID: ', 13376)
('In Thread :', 1, ', Process ID: ', 13375)
('In Thread :', 3, ', Process ID: ', 13377)
```


What can you summarize about the results? How does a 4 CPU machine compare to an 8 CPU machine in this exercise? [2 points]

Just like in case of 4 CPU machine, when the number of threads is less than or equal to the number of CPUs, then all threads get a full CPU and are able to utilize those CPUs fully. So in our case, when the number of threads were 2 and 4, then 2 and 4 CPUs were fully utilized whereas the rest were not utilized at all. Also, when we ran the program with 16 threads, all 8 CPUs were fully utilized and there was processes waiting for CPU time due to overloading.

4CPU vs 8 CPU

I think the main difference would be when we ran the program with 16 threads, in case of 4 CPU machine, only 4 threads would run at any time and the rest of 12 threads would be waiting for CPU time. Whereas in case of a 8 CPU machine, 8 threads would be running and only 8 threads would be waiting for CPU time. So, it would speed up the processing.

When we ran the program with 2 and 4 threads, both machines had equal or more CPUs than the number of threads and there were no other programs running on those machines, so all threads got full CPUs and I don't think the processing would have been any faster on the 8 CPU compared to 4 CPU machine.

Problem 2: I/O Analysis [points: 25]

Paste your source code into the following area. Make sure you clarify what you did to programmatically create an I/O intensive process. [10 points]

```
# This program creates a IO-Intensive workload.
# It calculates the Fibonacci series and writes a lot of text to files on disk

import argparse
import os
import time
import multiprocessing

#Using argparse library to parse the input arguments.
prog = "hw2"
desc = "Starts a specified number of threads that do CPU-Intesive work"
parser = argparse.ArgumentParser(prog=prog, description=desc)
parser.add_argument('--numThreads', '-n', default=4, type=int)

parsed_args = parser.parse_args()
numThreads = parsed_args.numThreads

print("Number of Threads: ", numThreads)

#Function that calculates the Fibonacci series
def Fibonacci_to_file(n):
    file_to_process = open("output.txt", "w")
    file_to_process.write("Inside fibonacci function")
```

```

file_to_process.close()
if n==1:
    return 0
elif n==2:
    return 1
else:
    return Fibonacci_to_file(n-1)+Fibonacci_to_file(n-2)

#Function that loops infinitely and puts load on the CPU.
def create_cpu_load(thread_num):
    while(True):
        time.sleep(1)
        print("In Thread :", thread_num, ", Process ID: ", os.getpid())
        Fibonacci_to_file(100)
        print("In Thread:", thread_num," ,Completed Fibonacci")

#The following code creates the number of threads specified in the input argument
threads= []
for thread_count_id in range(numThreads):
    t = multiprocessing.Process(target = create_cpu_load, args=(thread_count_id,))
    threads.append(t)
    t.start()

for current_thread in threads:
    current_thread.join()

```

Provide your table or graphs demonstrating the results of running this code with 2, 4 and 16 threads on a 4 CPU machine: [5 points]

```

[centos@ip-172-31-13-96 ~]$ python io_intensive_multi_processing.py --numThreads 2
('Number of Threads: ', 2)
('In Thread :', 0, ', Process ID: ', 17179)
('In Thread :', 1, ', Process ID: ', 17180)

```

TID	PRIO	USER	DISK READ	DISK WRITE	SWAPIN	IO>	COMMAND
17180	be/4	centos	0.00 B/s	2.61 M/s	0.00 %	52.93 %	python io_intensive_multi_processing.py --numThreads 2
17179	be/4	centos	0.00 B/s	4.20 M/s	0.00 %	36.38 %	python io_intensive_multi_processing.py --numThreads 2

```

1  [|||||] 5.5% Tasks: 33, 57 thr, 79 kthr; 1 running
2  [|||||] 12.8% Load average: 1.62 0.67 1.12
3  [|||||] 2.7% Uptime: 01:59:23
4  [|||||] 0.0%
Mem[|||||] 288M/7.15G
Swp[|||||] 0K/0K

```

PPID	CPU	TGID	PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
18839	3	18877	18877	centos	20	0	152M	7040	2268	S	0.0	0.1	0:00.01	python io_intensive_multi_processing.py --numThreads 2
18877	2	18879	18879	centos	20	0	152M	5576	744	D	11.5	0.1	0:13.55	python io_intensive_multi_processing.py --numThreads 2
18877	4	18878	18878	centos	20	0	152M	5580	748	D	13.5	0.1	0:14.52	python io_intensive_multi_processing.py --numThreads 2
1	1	858	858	root	20	0	549M	16832	6024	S	0.0	0.2	0:00.69	/usr/bin/python -Es /usr/sbin/tuned -l -P

```

[centos@ip-172-31-13-96 ~]$ python io_intensive_multi_processing.py --numThreads 4
('Number of Threads: ', 4)
('In Thread :', 0, ', Process ID: ', 17523)
('In Thread :', 2, ', Process ID: ', 17525)
('In Thread :', 1, ', Process ID: ', 17524)
('In Thread :', 3, ', Process ID: ', 17526)

```

```
hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 156x32
Total DISK READ :      0.00 B/s | Total DISK WRITE :      9.98 M/s
Actual DISK READ:      0.00 B/s | Actual DISK WRITE:      7.00 M/s
TID  PRIO  USER      DISK READ  DISK WRITE  SWAPIN     IO>     COMMAND
17524 be/4  centos    0.00 B/s   1596.18 K/s  0.00 %    23.30 % python io_intensive_multi_processing.py --numThreads 4
17523 be/4  centos    0.00 B/s   1897.20 K/s  0.00 %    21.27 % python io_intensive_multi_processing.py --numThreads 4
17525 be/4  centos    0.00 B/s   1691.24 K/s  0.00 %    20.92 % python io_intensive_multi_processing.py --numThreads 4
17526 be/4  centos    0.00 B/s   1738.77 K/s  0.00 %    20.16 % python io_intensive_multi_processing.py --numThreads 4
```

```
hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 153x52
1 [|||||] 4.9% Tasks: 35, 57 thr, 79 kthr; 2 running
2 [|||||] 11.5% Load average: 3.50 1.63 1.41
3 [|||||] 6.1% Uptime: 02:01:29
4 [|||||] 8.2%
Mem[|||||] 289M/7.15G
Swp[|||||] 0K/0K
PPID CPU TID PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
18839 3 18985 18985 centos 20 0 152M 7044 2268 S 0.0 0.1 0:00.01 python io_intensive_multi_processing.py --numThreads 4
18985 4 18989 18989 centos 20 0 152M 5584 744 D 10.2 0.1 0:10.86 python io_intensive_multi_processing.py --numThreads 4
18985 2 18988 18988 centos 20 0 152M 5584 744 D 10.2 0.1 0:10.66 python io_intensive_multi_processing.py --numThreads 4
18985 3 18987 18987 centos 20 0 152M 5580 744 D 11.6 0.1 0:10.99 python io_intensive_multi_processing.py --numThreads 4
18985 2 18986 18986 centos 20 0 152M 5584 748 D 10.2 0.1 0:10.90 python io_intensive_multi_processing.py --numThreads 4
1 3 858 858 root 20 0 549M 16832 6024 S 0.0 0.2 0:00.70 /usr/bin/python -Es /usr/sbin/tuned -l -P
```

```
[centos@ip-172-31-13-96 ~]$ python io_intensive_multi_processing.py --numThreads 16
('Number of Threads: ', 16)
('In Thread :', 0, 'Process ID: ', 17657)
('In Thread :', 2, 'Process ID: ', 17659)
('In Thread :', 1, 'Process ID: ', 17658)
('In Thread :', 3, 'Process ID: ', 17660)
('In Thread :', 4, 'Process ID: ', 17661)
('In Thread :', 5, 'Process ID: ', 17662)
('In Thread :', 6, 'Process ID: ', 17663)
('In Thread :', 7, 'Process ID: ', 17664)
('In Thread :', 8, 'Process ID: ', 17665)
('In Thread :', 10, 'Process ID: ', 17667)
('In Thread :', 11, 'Process ID: ', 17668)
('In Thread :', 9, 'Process ID: ', 17666)
('In Thread :', 13, 'Process ID: ', 17670)
('In Thread :', 12, 'Process ID: ', 17669)
('In Thread :', 15, 'Process ID: ', 17672)
('In Thread :', 14, 'Process ID: ', 17671)
```

```
hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 156x32
Total DISK READ :      0.00 B/s | Total DISK WRITE :      8.96 M/s
Actual DISK READ:      0.00 B/s | Actual DISK WRITE:      6.59 M/s
TID  PRIO  USER      DISK READ  DISK WRITE  SWAPIN     IO>     COMMAND
17669 be/4  centos    0.00 B/s   418.98 K/s  0.00 %    6.11 % python io_intensive_multi_processing.py --numThreads 16
17660 be/4  centos    0.00 B/s   371.55 K/s  0.00 %    6.11 % python io_intensive_multi_processing.py --numThreads 16
17666 be/4  centos    0.00 B/s   415.03 K/s  0.00 %    6.04 % python io_intensive_multi_processing.py --numThreads 16
17659 be/4  centos    0.00 B/s   466.41 K/s  0.00 %    5.87 % python io_intensive_multi_processing.py --numThreads 16
17658 be/4  centos    0.00 B/s   450.60 K/s  0.00 %    5.81 % python io_intensive_multi_processing.py --numThreads 16
17665 be/4  centos    0.00 B/s   411.07 K/s  0.00 %    5.68 % python io_intensive_multi_processing.py --numThreads 16
17670 be/4  centos    0.00 B/s   446.65 K/s  0.00 %    5.59 % python io_intensive_multi_processing.py --numThreads 16
17671 be/4  centos    0.00 B/s   395.26 K/s  0.00 %    5.58 % python io_intensive_multi_processing.py --numThreads 16
17663 be/4  centos    0.00 B/s   442.69 K/s  0.00 %    5.38 % python io_intensive_multi_processing.py --numThreads 16
17661 be/4  centos    0.00 B/s   438.74 K/s  0.00 %    5.09 % python io_intensive_multi_processing.py --numThreads 16
17672 be/4  centos    0.00 B/s   415.03 K/s  0.00 %    5.04 % python io_intensive_multi_processing.py --numThreads 16
17662 be/4  centos    0.00 B/s   395.26 K/s  0.00 %    4.83 % python io_intensive_multi_processing.py --numThreads 16
17657 be/4  centos    0.00 B/s   383.41 K/s  0.00 %    4.74 % python io_intensive_multi_processing.py --numThreads 16
17664 be/4  centos    0.00 B/s   375.50 K/s  0.00 %    4.57 % python io_intensive_multi_processing.py --numThreads 16
17667 be/4  centos    0.00 B/s   399.22 K/s  0.00 %    4.30 % python io_intensive_multi_processing.py --numThreads 16
17668 be/4  centos    0.00 B/s   339.93 K/s  0.00 %    3.65 % python io_intensive_multi_processing.py --numThreads 16
```



```

hw2 — centos@ip-172-31-13-96:~ — ssh -i kiran_e88.pem centos@ec2-18-220-128-51.us-east-2.compute.amazonaws.com — 153x52

 1 [|||||] 5.5% Tasks: 47, 57 thr, 79 kthr; 2 running
 2 [|||||] 11.6% Load average: 14.69 6.93 3.42
 3 [|||||] 6.1% Uptime: 02:04:04
 4 [|||||] 8.2%
Mem[|||||] 300M/7.15G
Swp[|||||] 0K/0K

PPID CPU TID PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
18839 3 19086 19086 centos 20 0 152M 7044 2268 S 0.0 0.1 0:00.02 python io_intensive_multi_processing.py --numThreads 16
19086 1 19102 19102 centos 20 0 152M 5604 744 D 2.7 0.1 0:03.59 python io_intensive_multi_processing.py --numThreads 16
19086 4 19101 19101 centos 20 0 152M 5600 744 D 2.7 0.1 0:03.63 python io_intensive_multi_processing.py --numThreads 16
19086 2 19100 19100 centos 20 0 152M 5600 744 D 2.0 0.1 0:03.59 python io_intensive_multi_processing.py --numThreads 16
19086 1 19099 19099 centos 20 0 152M 5596 744 D 2.7 0.1 0:03.57 python io_intensive_multi_processing.py --numThreads 16
19086 4 19098 19098 centos 20 0 152M 5596 744 D 2.7 0.1 0:03.63 python io_intensive_multi_processing.py --numThreads 16
19086 3 19097 19097 centos 20 0 152M 5596 744 D 2.7 0.1 0:03.60 python io_intensive_multi_processing.py --numThreads 16
19086 2 19096 19096 centos 20 0 152M 5592 744 D 2.7 0.1 0:03.60 python io_intensive_multi_processing.py --numThreads 16
19086 1 19095 19095 centos 20 0 152M 5592 744 D 3.4 0.1 0:03.61 python io_intensive_multi_processing.py --numThreads 16
19086 3 19094 19094 centos 20 0 152M 5592 744 D 2.7 0.1 0:03.61 python io_intensive_multi_processing.py --numThreads 16
19086 2 19093 19093 centos 20 0 152M 5588 744 D 2.7 0.1 0:03.60 python io_intensive_multi_processing.py --numThreads 16
19086 3 19092 19092 centos 20 0 152M 5588 744 D 2.0 0.1 0:03.57 python io_intensive_multi_processing.py --numThreads 16
19086 4 19091 19091 centos 20 0 152M 5580 744 D 2.7 0.1 0:03.61 python io_intensive_multi_processing.py --numThreads 16
19086 3 19090 19090 centos 20 0 152M 5580 744 D 2.7 0.1 0:03.60 python io_intensive_multi_processing.py --numThreads 16
19086 2 19089 19089 centos 20 0 152M 5580 744 D 2.7 0.1 0:03.60 python io_intensive_multi_processing.py --numThreads 16
19086 4 19088 19088 centos 20 0 152M 5576 744 D 2.7 0.1 0:03.57 python io_intensive_multi_processing.py --numThreads 16
19086 2 19087 19087 centos 20 0 152M 5580 748 D 3.4 0.1 0:03.60 python io_intensive_multi_processing.py --numThreads 16
1 3 858 858 root 20 0 549M 16832 6024 S 0.0 0.2 0:00.72 /usr/bin/python -Es /usr/sbin/tuned -l -P

```

Number of Threads	Actual DISK Write	Load over last min
2	7.11 M/s	1.62
4	7.00 M/s	3.50
16	6.59 M/s	14.69

What can you summarize about the results? [3 points]

It appears that the actual disk writes is decreasing as the number of threads increases. Also CPUs are not fully utilized even when number of threads is greater than or equal to number of CPUs.

Provide your table or graphs demonstrating the results of running this code with 2, 4 and 16 threads on an 8 CPU machine: [5 points]

```

[centos@ip-172-31-13-250 ~]$ python io_intensive_multi_processing.py --numThreads 2
('Number of Threads: ', 2)
('In Thread :', 0, ', Process ID: ', 15574)
('In Thread :', 1, ', Process ID: ', 15575)

hw2 — centos@ip-172-31-13-250:~ — ssh -i kiran_e88.pem centos@ec2-13-59-183-19.us-east-2.compute.amazonaws.com — 156x32
Total DISK READ : 0.00 B/s | Total DISK WRITE : 13.72 M/s
Actual DISK READ: 0.00 B/s | Actual DISK WRITE: 7.18 M/s
TID PRIO USER DISK READ DISK WRITE SWAPIN IO> COMMAND
15575 be/4 centos 0.00 B/s 3.07 M/s 0.00 % 47.97 % python io_intensive_multi_processing.py --numThreads 2
15574 be/4 centos 0.00 B/s 3.81 M/s 0.00 % 40.79 % python io_intensive_multi_processing.py --numThreads 2

```

```
hw2 — centos@ip-172-31-13-250:~ — ssh -i kiran_e88.pem centos@ec2-13-59-183-19.us-east-2.compute.amazonaws.com — 153x52

1 [|||||] 1.4% Tasks: 38, 70 thr, 105 kthr; 2 running
2 [|||||] 8.2% Load average: 2.34 6.19 6.60
3 [|||||] 7.5% Uptime: 01:12:55
4 [|||||] 6.9%
5 [|||||] 1.3%
6 [|||||] 0.0%
7 [|||||] 0.0%
8 [|||||] 0.0%
Mem[|||||] 355M/14.5G
Swp[|||||] 0K/0K

PPID CPU TGID PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
15509 5 16450 16450 centos 20 0 152M 7044 2268 S 0.0 0.0 0:00.01 python io_intensive_multi_processing.py --numThreads 2
16450 3 16452 16452 centos 20 0 152M 5580 744 D 17.6 0.0 0:17.37 python io_intensive_multi_processing.py --numThreads 2
16450 4 16451 16451 centos 20 0 152M 5584 748 D 18.2 0.0 0:18.15 python io_intensive_multi_processing.py --numThreads 2
15560 1 15561 15561 root 20 0 193M 11832 4124 S 1.4 0.1 0:11.92 /usr/bin/python /sbin/iotop --only
1 5 896 896 root 20 0 549M 16652 5920 S 0.0 0.1 0:00.47 /usr/bin/python -Es /usr/sbin/tuned -l -P
```

```
[centos@ip-172-31-13-250 ~]$ python io_intensive_multi_processing.py --numThreads 4
('Number of Threads: ', 4)
('In Thread :', 0, ', Process ID: ', 15631)
('In Thread :', 2, ', Process ID: ', 15633)
('In Thread :', 1, ', Process ID: ', 15632)
('In Thread :', 3, ', Process ID: ', 15634)
```

```
hw2 — centos@ip-172-31-13-250:~ — ssh -i kiran_e88.pem centos@ec2-13-59-183-19.us-east-2.compute.amazonaws.com — 156x32

Total DISK READ : 0.00 B/s | Total DISK WRITE : 11.70 M/s
Actual DISK READ: 0.00 B/s | Actual DISK WRITE: 6.53 M/s

TID PRIO USER DISK READ DISK WRITE SWAPIN IO% COMMAND
15631 be/4 centos 0.00 B/s 1590.06 K/s 0.00 % 22.01 % python io_intensive_multi_processing.py --numThreads 4
15634 be/4 centos 0.00 B/s 1491.17 K/s 0.00 % 21.99 % python io_intensive_multi_processing.py --numThreads 4
15632 be/4 centos 0.00 B/s 1760.14 K/s 0.00 % 20.30 % python io_intensive_multi_processing.py --numThreads 4
15633 be/4 centos 0.00 B/s 1641.48 K/s 0.00 % 18.85 % python io_intensive_multi_processing.py --numThreads 4
```

```
hw2 — centos@ip-172-31-13-250:~ — ssh -i kiran_e88.pem centos@ec2-13-59-183-19.us-east-2.compute.amazonaws.com — 153x52

1 [|||||] 4.8% Tasks: 40, 70 thr, 105 kthr; 1 running
2 [|||||] 10.0% Load average: 4.22 8.36 7.25
3 [|||||] 8.8% Uptime: 01:10:45
4 [|||||] 10.0%
5 [|||||] 2.0%
6 [|||||] 0.0%
7 [|||||] 1.3%
8 [|||||] 1.4%
Mem[|||||] 356M/14.5G
Swp[|||||] 0K/0K

PPID CPU TGID PID USER PRI NI VIRT RES SHR S CPU% MEM% TIME+ Command
15509 1 16242 16242 centos 20 0 152M 7040 2268 S 0.0 0.0 0:00.01 python io_intensive_multi_processing.py --numThreads 4
16242 2 16246 16246 centos 20 0 152M 5580 744 D 12.1 0.0 0:25.22 python io_intensive_multi_processing.py --numThreads 4
16242 7 16245 16245 centos 20 0 152M 5580 744 D 12.1 0.0 0:25.22 python io_intensive_multi_processing.py --numThreads 4
16242 1 16244 16244 centos 20 0 152M 5576 744 D 11.4 0.0 0:25.36 python io_intensive_multi_processing.py --numThreads 4
16242 3 16243 16243 centos 20 0 152M 5580 748 D 12.1 0.0 0:25.16 python io_intensive_multi_processing.py --numThreads 4
15560 5 15561 15561 root 20 0 193M 11832 4124 S 1.3 0.1 0:10.62 /usr/bin/python /sbin/iotop --only
1 5 896 896 root 20 0 549M 16652 5920 S 0.0 0.1 0:00.45 /usr/bin/python -Es /usr/sbin/tuned -l -P
```

```
[centos@ip-172-31-13-250 ~]$ python io_intensive_multi_processing.py --numThreads 16
('Number of Threads: ', 16)
('In Thread :', 0, ', Process ID: ', 15719)
('In Thread :', 2, ', Process ID: ', 15721)
('In Thread :', 1, ', Process ID: ', 15720)
('In Thread :', 4, ', Process ID: ', 15723)
('In Thread :', 3, ', Process ID: ', 15722)
('In Thread :', 6, ', Process ID: ', 15725)
('In Thread :', 5, ', Process ID: ', 15724)
('In Thread :', 8, ', Process ID: ', 15727)
('In Thread :', 7, ', Process ID: ', 15726)
('In Thread :', 10, ', Process ID: ', 15729)
('In Thread :', 9, ', Process ID: ', 15728)
('In Thread :', 12, ', Process ID: ', 15731)
('In Thread :', 11, ', Process ID: ', 15730)
('In Thread :', 14, ', Process ID: ', 15733)
('In Thread :', 13, ', Process ID: ', 15732)
('In Thread :', 15, ', Process ID: ', 15734)
```

```

hw2 — centos@ip-172-31-13-250:~ — ssh -i kiran_e88.pem centos@ec2-13-59-183-19.us-east-2.compute.amazonaws.com — 156x32
Total DISK READ : 0.00 B/s | Total DISK WRITE : 10.76 M/s
Actual DISK READ: 0.00 B/s | Actual DISK WRITE: 6.40 M/s

```

TID	PRI	USER	DISK READ	DISK WRITE	SWAPIN	IO>	COMMAND
15719	be/4	centos	0.00 B/s	367.28 K/s	0.00 %	6.33 %	python io_intensive_multi_processing.py --numThreads 16
15728	be/4	centos	0.00 B/s	387.03 K/s	0.00 %	6.10 %	python io_intensive_multi_processing.py --numThreads 16
15726	be/4	centos	0.00 B/s	418.62 K/s	0.00 %	5.59 %	python io_intensive_multi_processing.py --numThreads 16
15727	be/4	centos	0.00 B/s	351.49 K/s	0.00 %	5.57 %	python io_intensive_multi_processing.py --numThreads 16
15729	be/4	centos	0.00 B/s	355.44 K/s	0.00 %	5.53 %	python io_intensive_multi_processing.py --numThreads 16
15730	be/4	centos	0.00 B/s	422.57 K/s	0.00 %	5.44 %	python io_intensive_multi_processing.py --numThreads 16
15732	be/4	centos	0.00 B/s	355.44 K/s	0.00 %	5.42 %	python io_intensive_multi_processing.py --numThreads 16
15723	be/4	centos	0.00 B/s	383.08 K/s	0.00 %	5.31 %	python io_intensive_multi_processing.py --numThreads 16
15734	be/4	centos	0.00 B/s	383.08 K/s	0.00 %	5.16 %	python io_intensive_multi_processing.py --numThreads 16
15722	be/4	centos	0.00 B/s	430.47 K/s	0.00 %	5.07 %	python io_intensive_multi_processing.py --numThreads 16
15725	be/4	centos	0.00 B/s	406.78 K/s	0.00 %	5.04 %	python io_intensive_multi_processing.py --numThreads 16
15724	be/4	centos	0.00 B/s	434.42 K/s	0.00 %	4.86 %	python io_intensive_multi_processing.py --numThreads 16
15721	be/4	centos	0.00 B/s	410.73 K/s	0.00 %	4.82 %	python io_intensive_multi_processing.py --numThreads 16
15731	be/4	centos	0.00 B/s	489.71 K/s	0.00 %	4.60 %	python io_intensive_multi_processing.py --numThreads 16
15733	be/4	centos	0.00 B/s	398.88 K/s	0.00 %	4.43 %	python io_intensive_multi_processing.py --numThreads 16
15720	be/4	centos	0.00 B/s	375.18 K/s	0.00 %	4.19 %	python io_intensive_multi_processing.py --numThreads 16

```

hw2 — centos@ip-172-31-13-250:~ — ssh -i kiran_e88.pem centos@ec2-13-59-183-19.us-east-2.compute.amazonaws.com — 153x52

```

```

1 [|||||] 4.8% Tasks: 52, 70 thr, 103 kthr; 3 running
2 [|||||] 9.6% Load average: 16.14 13.41 7.90
3 [|||||] 8.1% Uptime: 01:06:07
4 [|||||] 8.8%
5 [|||||] 1.4%
6 [|||||] 2.0%
7 [|||||] 2.0%
8 [|||||] 0.7%
Mem[|||||] 369M/14.5G
Swp[|||||] 0K/0K

```

PPID	CPU	TGID	PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
15509	5	15718	15718	centos	20	0	152M	7044	2268	S	0.0	0.0	0:00.02	python io_intensive_multi_processing.py --numThreads 16
15718	6	15734	15734	centos	20	0	152M	5604	744	D	3.4	0.0	0:15.37	python io_intensive_multi_processing.py --numThreads 16
15718	2	15733	15733	centos	20	0	152M	5600	744	D	3.4	0.0	0:15.35	python io_intensive_multi_processing.py --numThreads 16
15718	6	15732	15732	centos	20	0	152M	5600	744	D	3.4	0.0	0:15.49	python io_intensive_multi_processing.py --numThreads 16
15718	2	15731	15731	centos	20	0	152M	5596	744	D	2.7	0.0	0:15.36	python io_intensive_multi_processing.py --numThreads 16
15718	8	15730	15730	centos	20	0	152M	5596	744	D	2.7	0.0	0:15.29	python io_intensive_multi_processing.py --numThreads 16
15718	5	15729	15729	centos	20	0	152M	5596	744	D	3.4	0.0	0:15.69	python io_intensive_multi_processing.py --numThreads 16
15718	3	15728	15728	centos	20	0	152M	5592	744	D	2.7	0.0	0:15.51	python io_intensive_multi_processing.py --numThreads 16
15718	7	15727	15727	centos	20	0	152M	5592	744	D	3.4	0.0	0:15.49	python io_intensive_multi_processing.py --numThreads 16
15718	1	15726	15726	centos	20	0	152M	5592	744	D	2.7	0.0	0:15.36	python io_intensive_multi_processing.py --numThreads 16
15718	3	15725	15725	centos	20	0	152M	5588	744	D	2.7	0.0	0:15.66	python io_intensive_multi_processing.py --numThreads 16
15718	2	15724	15724	centos	20	0	152M	5584	744	D	2.0	0.0	0:15.46	python io_intensive_multi_processing.py --numThreads 16
15718	2	15723	15723	centos	20	0	152M	5580	744	D	3.4	0.0	0:15.41	python io_intensive_multi_processing.py --numThreads 16
15718	8	15722	15722	centos	20	0	152M	5580	744	D	3.4	0.0	0:15.24	python io_intensive_multi_processing.py --numThreads 16
15718	3	15721	15721	centos	20	0	152M	5580	744	R	3.4	0.0	0:15.34	python io_intensive_multi_processing.py --numThreads 16
15718	1	15720	15720	centos	20	0	152M	5576	744	D	4.1	0.0	0:15.56	python io_intensive_multi_processing.py --numThreads 16
15718	2	15719	15719	centos	20	0	152M	5580	748	R	2.7	0.0	0:15.31	python io_intensive_multi_processing.py --numThreads 16
15560	4	15561	15561	root	20	0	193M	11832	4124	S	1.4	0.1	0:07.68	/usr/bin/python /sbin/ioping --only

Number of Threads	Actual DISK write	Load over last minute
2	7.18 M/s	2.34
4	6.53 M/s	4.22
16	6.40 M/s	16.14

What can you summarize about the results? How does a 4 CPU machine compare to an 8 CPU machine in this exercise? [2 points]

The results from 4 CPU machine are very similar to 8 CPU machines and it appears we have reached the IO limit in both cases. Also, for both machines as the number of threads increased the actual disk write has decreased.

Problem 3: unique counts [points: 25]

Paste your source code into the following area [10 points]

```
import csv
from dateutil.parser import parse
import time
from multiprocessing import Process, Lock, Manager

files = ['file-input1.csv', 'file-input2.csv', 'file-input3.csv', 'file-input4.csv']

# Following function reads each line in the above files and forms a data structure of all data
def readFile(inputFile, lock, shared_dict):
    lock.acquire()
    with open(inputFile) as csvfile:
        readCSV = csv.reader(csvfile, delimiter=',')
        for row in readCSV:
            dt = parse(row[1])
            dt_str = str(dt.date())+":"+str(dt.time().hour)
            url = row[2]
            user = row[3]
            if dt_str in shared_dict:
                if url in shared_dict[dt_str]:
                    if user in shared_dict[dt_str][url]:
                        dictCopy = shared_dict[dt_str][url]
                        dictCopy[user] = dictCopy[user]+1
                        dictSecondCopy = shared_dict[dt_str]
                        dictSecondCopy[url] = dictCopy
                        shared_dict[dt_str] = dictSecondCopy
                        shared_dict[dt_str][url] = dictCopy
                    else:
                        ### First modify at the user level.
                        dictCopy = shared_dict[dt_str][url]
                        dictCopy[user] = 1
                        ### Next modify at the URL level
                        dictSecondCopy = shared_dict[dt_str]
                        dictSecondCopy[url] = dictCopy
                        shared_dict[dt_str] = dictSecondCopy
                        shared_dict[dt_str][url] = dictCopy
                else:
                    dictCopy = shared_dict[dt_str]
                    dictCopy[url] = {user: 1}
                    shared_dict[dt_str] = dictCopy
            else:
                shared_dict[dt_str] = {url: {user: 1}}
    lock.release()

manager = Manager()
shared_dict = manager.dict()
lock = Lock()

threads = []
for inputFile in files:
    procs = Process(target=readFile, args=[inputFile, lock, shared_dict])
    threads.append(procs)
    procs.start()

for t in threads : t.join()

####Query 1 #####
print("##### Query 1 OUTPUT #####")
for key in sorted(shared_dict.keys()):
    print (key, len(shared_dict[key]))
```



```

####Query 2 #####

print("##### Query 2 OUTPUT #####")
for key in sorted(shared_dict.keys()):
    for url in shared_dict[key]:
        print (key, url , len(shared_dict[key][url]))

####Query 3 #####

print("##### Query 3 OUTPUT #####")

for key in sorted(shared_dict.keys()):
    for url in shared_dict[key]:
        counter = 0
        for user in shared_dict[key][url]:
            counter = counter + shared_dict[key][url][user]
        print (key, url , counter)

```

Explain your choice of the data structures for shared state management [5 points]

I chose a nested dictionary with the following format:
 {date : { URL : { user : clicks}}}

Using this data structure allows us to answer all the following queries using the same program/data structure in an efficient manner.

What are the results of your queries for the following specified keys? [10 points]
 The expected output for the first value is provided for your reference.

Query 1:

<date_hour>, <url_count>

```

2019-09-12:13, 185
2019-09-12:14, 186
2019-09-12:15, 185
2019-09-12:16, 190
2019-09-12:17, 189

```

Query 2

<date:hour:url>, unique_user_count

```

2019-09-12:02:http://example.com/?url=003, 1
2019-09-12:02:http://example.com/?url=004, 3
2019-09-12:02:http://example.com/?url=005, 4
2019-09-12:02:http://example.com/?url=006, 10

```

Query 3

<date:hour:url>, event_count

```

2019-09-12:02:http://example.com/?url=003, 1
2019-09-12:02:http://example.com/?url=004, 3
2019-09-12:02:http://example.com/?url=005, 5

```

```
2019-09-12:02:http://example.com/?url=006, 10
```

Problem 4: time range queries [points: 25]

Paste your source code into the following area [15 points]

```
import csv
from dateutil.parser import parse
from datetime import datetime
import time
from multiprocessing import Process, Lock, Manager

files = ['file-input1.csv', 'file-input2.csv', 'file-input3.csv', 'file-input4.csv']

# Following function reads each line in the above files and forms a data structure of all the data
def readFile(inputFile, lock, shared_dict):
    lock.acquire()
    with open(inputFile) as csvfile:
        readCSV = csv.reader(csvfile, delimiter=',')
        for row in readCSV:
            dt = parse(row[1])
            dt_str = str(dt.date())+" "+str(dt.time().hour)
            url = row[2]
            country = row[4]

            start_date = datetime(2019, 9, 13, 17)
            end_date = datetime(2019, 9, 14, 9)
            current_line_date = datetime(dt.year, dt.month, dt.day, dt.time().hour)
            if (start_date <= current_line_date <= end_date):
                if dt_str in shared_dict:
                    if country in shared_dict[dt_str]:
                        if url in shared_dict[dt_str][country]:
                            dictCopy = shared_dict[dt_str][country]
                            dictCopy[url] = dictCopy[url]+1
                            dictSecondCopy = shared_dict[dt_str]
                            dictSecondCopy[country] = dictCopy
                            shared_dict[dt_str] = dictSecondCopy
                            shared_dict[dt_str][country] = dictCopy
                        else:
                            ### First modify at the user level.
                            dictCopy = shared_dict[dt_str][country]
                            dictCopy[url] = 1
                            ### Next modify at the URL level
                            dictSecondCopy = shared_dict[dt_str]
                            dictSecondCopy[country] = dictCopy
                            shared_dict[dt_str] = dictSecondCopy
                            shared_dict[dt_str][country] = dictCopy
                    else:
                        dictCopy = shared_dict[dt_str]
                        dictCopy[country] = {url: 1}
                        shared_dict[dt_str] = dictCopy
                else:
                    shared_dict[dt_str] = {country: {url: 1}}

    lock.release()

manager = Manager()
shared_dict = manager.dict()
lock = Lock()
```

```

threads = []
for inputFile in files:
    procs = Process(target=readFile, args=[inputFile, lock, shared_dict])
    threads.append(procs)
    procs.start()

for t in threads : t.join()

####Query 1 #####
print("##### Query 1 OUTPUT #####")
for key in sorted(shared_dict.keys()):
    for country in shared_dict[key]:
        print (key, country, len(shared_dict[key][country]))

```

What are the main differences with the Problem 3 implementation? [5 points]

I modified the data structure as follows:
 {date: {country: {url: count}}}

Also, I have to filter out the data to make sure we only get the time range we want. This is done before adding the data to dictionary.

What are the results of your query for the specified keys ? [5 points]

The expected output for the first value is provided for your reference.

```

<date,hour,country>, url_count
2019-09-13:19,IQ, 1
2019-09-13:19,IR, 4
2019-09-13:19,IS, 9
2019-09-13:19,IT, 2
2019-09-13:19,JE, 4

```

Problem 5: Bonus: Top N queries [15 points]

Paste your source code into the following area [5 points]

What are the main differences with the Problem 3 and 4 implementation? [5 points]

What are the results of your query? [5 points] The expected 5 values for 9/12 are provided, please fill in the values for avg TTFB and the URLs for 9/13 and 9/14.

Date	URL	Average_TTFB
9/12/19	http://example.com/?url=114	0.393101408

9/12/19	http://example.com/?url=101	0.402545
9/12/19	http://example.com/?url=133	0.413317187
9/12/19	http://example.com/?url=033	0.418867857
9/12/19	http://example.com/?url=157	0.419289394

9/13/19
9/13/19
9/13/19
9/13/19
9/13/19

9/14/19
9/14/19
9/14/19
9/14/19
9/14/19